

UC Santa Barbara

UC Santa Barbara Electronic Theses and Dissertations

Title

Regularization and Look-Ahead Procedures for Selection of Basis Functions from Multiple Libraries

Permalink

<https://escholarship.org/uc/item/6ss3v7cb>

Author

Zhu, Ling

Publication Date

2017

Peer reviewed|Thesis/dissertation

University of California
Santa Barbara

Regularization and Look-Ahead Procedures for Selection of Basis Functions from Multiple Libraries

A dissertation submitted in partial satisfaction
of the requirements for the degree

Doctor of Philosophy
in
Statistics

by

Ling Zhu

Committee in charge:

Professor Yuedong Wang, Chair
Professor Wendy Meiring
Professor John Hsu

March 2017

The Dissertation of Ling Zhu is approved.

Professor Wendy Meiring

Professor John Hsu

Professor Yuedong Wang, Committee Chair

March 2017

Regularization and Look-Ahead Procedures for Selection of Basis Functions from
Multiple Libraries

Copyright © 2017

by

Ling Zhu

*I want to dedicate this work to my parents,
for their unwavering love and encouragement.*

Acknowledgements

First and foremost, I would like to express my most sincere gratitude to my advisor, Professor Yuedong Wang. His dedication to teaching and statistical research has constantly inspired me during my graduate study. I am very grateful for all the knowledge I gained from attending the excellent classes he has offered. His patient guidance and valuable insight have helped me conquer numerous obstacles in my research and made the work presented in dissertation possible. I would also like to thank my committee members, Professor Wendy Meiring and Professor John Hsu, who gave me great encouragement and advice. I am indebted to Professor Meiring for her tremendous help in improving the accuracy of my writing, as well as her expert opinions regarding the real application of my methodology.

My deep appreciation also goes to all other faculty members, staff, and colleagues in the Statistics department. Together you have created a nourishing and friendly academic environment that helped me grow. I will cherish every moment I spent in the research seminars and every conversation I had with my colleagues regarding research and life. We always had support for each other and thanks to you my time at UC Santa Barbara was absolutely wonderful.

Finally, I want to thank my parents, Jianxun Zhu and Hongmei Cai, for their unconditional love and support. Their unwavering encouragement gave me strength and confidence. I owe everything to them.

Curriculum Vitæ

Ling Zhu

Education

- 2016 Doctor of Philosophy in Statistics and Applied Probability, Department of Statistics and Applied Probability, University of California, Santa Barbara.
- 2010 Master of Arts in Economics, Department of Economics, University of California, Santa Barbara.
- 2009 Bachelor of Arts in Economics, Zhejiang Gongshang University, Hangzhou, China.

Education

- 2012-2016 Teaching Assistant, Department of Statistics and Applied Probability, University of California, Santa Barbara.
- 2013 Research Intern, EdLab, Teachers College, Columbia University, New York, New York.
- 2011 Readership Assistant, Department of Statistics and Applied Probability, University of California, Santa Barbara.

Abstract

Regularization and Look-Ahead Procedures for Selection of Basis Functions from
Multiple Libraries

by

Ling Zhu

Basis Selection from Multiple Libraries (BSML) procedure was proposed by Sklar et al. (2013) to estimate spatially inhomogeneous functions with linear combinations of basis functions that are adaptively selected from multiple libraries. This methodology proves to be successful but suffers from certain drawbacks. First, the BSML procedures utilize the generalized degrees of freedom (GDF) to measure the complexity of a modeling procedure. This approach can be computationally demanding, since the GDF needs to be estimated at every step of the forward selection process. Another drawback of the BSML procedures is its greedy nature when searching for basis functions. At each forward selection step, the BSML procedures search for only one basis function to add to the current model.

We first propose two procedures called Adaptive LASSO Basis Selection (ALBS) and its modified version ALBS-2 to address the first drawback. The idea is to penalize the bases differently based on their induced reductions of RSS and the libraries they are from. This regularization approach shows no clear advantages over the BSML procedures in terms of performance according to our simulations.

We then propose the look-ahead procedure (LAP) to address both drawbacks of the BSML procedures. LAP avoids the estimation of GDF by treating the inflated degrees of freedom (IDF) of each library as tuning parameters, and then estimate them using cross validation. Moreover, LAP adopts less greedy search rules in the forward selection

so that either one or a pair of basis functions can be added to the current model after each iteration. Extensive simulations show that the LAP can outperform the BSML procedures and other widely-used modeling procedures both in terms of computation speed and mean squared error. Interesting real data applications in geology and meteorology using the look-ahead procedure are also provided.

Contents

Curriculum Vitae	vi
Abstract	vii
1 Introduction	1
1.1 Nonparametric Regression with Splines	1
1.2 Smoothing Spline ANOVA	7
1.3 Component Selection and Smoothing Operator	9
1.4 LASSO and Adaptive LASSO	11
1.5 Hybrid Adaptive Spline	14
1.6 Multivariate Adaptive Regression Splines	16
1.7 Generalized Degrees of Freedom and Covariance Penalty	20
1.8 The BSML Procedure	25
2 Basis Selection from Multiple Libraries Using Adaptive LASSO	33
2.1 Adaptive LASSO Basis Selection	33
2.2 Simulations to Compare ALBS, BSML, and HAS	36
2.3 Adaptive LASSO Basis Selection with Estimated IDF	41
2.4 Simulations to Compare ALBS-2, BSML, and HAS	42
3 Look-Ahead Procedure	46
3.1 Problems With Greedy Search	46
3.2 Forward Selection In the Look-Ahead Procedure with Fixed IDFs	48
3.3 Forward Selection Via Householder Transformation	52
3.4 Forward Selection In the Look-Ahead Procedure with IDFs as Tuning Parameters	70
3.5 Selection of IDFs	82
3.6 The Whole Look-Ahead Procedure	84
3.7 Bootstrap Confidence Intervals	87
3.8 R Functions for the Look-Ahead Procedure	91

4	Simulations	98
4.1	Estimation of Univariate Functions	98
4.2	Bootstrap Confidence Intervals	112
4.3	Estimation of Multivariate Functions	118
4.4	Three Other Examples	131
5	Applications	140
5.1	Well Log	140
5.2	Ozone Pressure	146

Chapter 1

Introduction

1.1 Nonparametric Regression with Splines

Consider the following univariate nonparametric regression problem:

$$y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n, \tag{1.1}$$

where y_i and x_i are the i th observation of the response variable Y and independent variable X respectively, ε_i 's are random errors with mean zero and variance σ^2 , and f is the regression function to be estimated on the domain $\mathcal{X} = [0, 1]$. We do not assume f to have a predetermined form. Instead, we only require f to have certain qualitative properties, such as smoothness, and let data decide the form of f . Thus, the model considered in (1.1) is nonparametric. There are two basic spline methods that can be used to estimate f : regression spline and smoothing spline.

1.1.1 Regression Spline

The regression spline model assumes that f can be well-approximated by either polynomial spline, or natural polynomial spline. Let $\pi^r = \{f : f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_{r-1} x^{r-1}\}$ be the collection of polynomials of order r , and C^r be the set of functions such that their r -th derivatives are continuous at each x . Then a polynomial spline f of order r with interior knots $\mathbf{t} = (t_1, \dots, t_k)$, where $0 = t_0 < t_1 < \cdots < t_k < t_{k+1} = 1$, satisfies:

- (i) $f(x) \in \pi^r$, for $x \in [t_i, t_{i+1}]$, $i = 0, \dots, k$.
- (ii) $f \in C^{r-2}$, and its $(r-1)$ -st derivative is a step function with jumps at the knots.

We denote the collection of all such polynomial splines of order r with interior knots \mathbf{t} as $S^r(\mathbf{t})$. For f to be a natural polynomial spline, it has to satisfy an additional condition:

- (iii) When $r = 2m$, $f(x) \in \pi^m$ for $x \in [0, t_1]$ and $x \in [t_k, 1]$.

This translates to the boundary conditions on the natural splines: $f^{(j)}(0) = f^{(j)}(1) = 0$, for $j = m, \dots, 2m-1$, where $f^{(j)}$ is the j th derivative of f . We denote the collection of all such natural splines of order $2m$ with interior knots \mathbf{t} as $NS^{2m}(\mathbf{t})$.

The space $S^r(\mathbf{t})$ is spanned by $r+k$ basis functions: $\{1, x, \dots, x^{r-1}, (x-t_1)_+^{r-1}, \dots, (x-t_k)_+^{r-1}\}$, so for any $f \in S^r(\mathbf{t})$, it can be uniquely represented by

$$f(x) = \sum_{j=0}^{r-1} \theta_j x^j + \sum_{j=1}^k \delta_j (x - t_j)_+^{r-1}. \quad (1.2)$$

Let $\mathbf{X}_{\mathbf{t}}$ be the $n \times (r+k)$ design matrix with the basis functions of $S^r(\mathbf{t})$ as its columns, where the i th row corresponds to the $r+k$ basis functions evaluated at observation location x_i for $i = 1, \dots, n$. Then (1.1) can be written as:

$$\mathbf{y} = \mathbf{X}_{\mathbf{t}} \boldsymbol{\beta}_{\mathbf{t}} + \boldsymbol{\varepsilon}, \quad (1.3)$$

where $\mathbf{y} = (y_1, \dots, y_n)'$, $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$, and $\boldsymbol{\beta}_t = (\theta_0, \dots, \theta_{r-1}, \delta_1, \dots, \delta_k)'$. The estimate of $\boldsymbol{\beta}_t$ that minimizes the residual sum of squares $\text{RSS}(\mathbf{t}) = (\mathbf{y} - \mathbf{X}_t \boldsymbol{\beta}_t)'(\mathbf{y} - \mathbf{X}_t \boldsymbol{\beta}_t)$ is:

$$\hat{\boldsymbol{\beta}}_t = (\mathbf{X}_t' \mathbf{X}_t)^{-1} \mathbf{X}_t' \mathbf{y}. \quad (1.4)$$

The solution in (1.4) is for fixed knots $\mathbf{t} = (t_1, \dots, t_k)$. However, both the number and the locations of the knots in \mathbf{t} should be treated as tuning parameters. Some ad hoc rules are available for selecting \mathbf{t} adaptively (See Eubank (1999)). The generalized cross validation (GCV) criterion can also be used to select both the knot locations and number of interior knots (k):

$$\text{GCV}(\mathbf{t}) = \frac{\text{RSS}(\mathbf{t})}{(n - (r + k))^2}, \quad (1.5)$$

where $r + k$ is the total number of basis functions. The tuning parameter $\hat{\mathbf{t}}$ is selected to minimize the GCV in (1.5).

The scenario is similar when we assume $f \in NS^{2m}(\mathbf{t})$. It is obvious that $NS^{2m}(\mathbf{t}) \subset S^{2m}(\mathbf{t})$ since a natural spline has to satisfy the boundary condition (iii) as well. Because there are a total of $2m$ such boundary constraints, the dimension of $NS^{2m}(\mathbf{t})$ is $r + k - 2m = k$, i.e., $NS^{2m}(\mathbf{t})$ is spanned by only k basis functions. Any $f \in NS^{2m}(\mathbf{t})$ can be represented by

$$f(x) = \sum_{j=0}^{m-1} \theta_j x^j + \sum_{j=1}^k \delta_j (x - t_j)_+^{2m-1}. \quad (1.6)$$

1.1.2 Smoothing Spline

The smoothing spline model assumes f to be in certain functional space, such as the Sobolev space $W_2^m[0, 1]$:

$$W_2^m[0, 1] = \left\{ f : f, f^{(1)}, \dots, f^{(m-1)} \text{ are absolutely continuous, } \int_0^1 (f^{(m)}(x))^2 dx < \infty \right\}. \quad (1.7)$$

Then f is estimated by solving the following penalized least square (PLS) problem:

$$\hat{f}_\lambda = \underset{f \in W_2^m[0,1]}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int_0^1 (f^{(m)}(x))^2 dx, \quad (1.8)$$

where λ is the tuning parameter that balances the trade-off between the goodness of fit and the smoothness of \hat{f}_λ . Assuming the design points are distinct, then for a fixed λ , the solution \hat{f}_λ belongs to $NS^{2m}(x_1, \dots, x_n)$ (See Eubank (1999)). When $\lambda = 0$, \hat{f}_λ is a natural spline that interpolates the data, since there is no penalty on the roughness of \hat{f}_λ . On the other hand, when $\lambda = \infty$, \hat{f}_λ becomes an m -th order polynomial. A commonly used criterion for selecting the smoothing parameter λ is the GCV (See Craven and Wahba (1978)). Let $\hat{\mathbf{f}}_\lambda = (\hat{f}_\lambda(x_1), \dots, \hat{f}_\lambda(x_n))'$, and \mathbf{P}_λ be the hat matrix that satisfies $\hat{\mathbf{f}}_\lambda = \mathbf{P}_\lambda \mathbf{y}$, then the GCV criterion is defined as

$$\text{GCV}(\lambda) = \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_\lambda(x_i))^2}{\left[\frac{1}{n} \operatorname{tr}(\mathbf{I} - \mathbf{P}_\lambda)\right]^2}, \quad (1.9)$$

where \mathbf{I} is the $n \times n$ identity matrix, and $\operatorname{tr}(\cdot)$ returns the trace of a matrix. The GCV estimate of λ is the minimizer of $\text{GCV}(\lambda)$. Other criteria such as the unbiased risk (UBR) criterion and generalized maximum likelihood (GML) criterion can also be used to select λ (see Wahba (1990)).

The Sobolev space $W_2^m[0,1]$ is an example of a reproducing kernel Hilbert space (RKHS), where a reproducing kernel (RK) $R(x, \cdot)$ is defined. More generally, let \mathcal{H} be a RKHS with the domain \mathcal{X} , where \mathcal{X} is an arbitrary set, then $R(x, \cdot)$ is in \mathcal{H} with the reproducing property:

$$\langle R(x, \cdot), f \rangle = f(x) \quad (1.10)$$

for any $x \in \mathcal{X}$ and $f \in \mathcal{H}$, where $\langle \cdot, \cdot \rangle$ is the inner product on \mathcal{H} . $R(x, \cdot)$ is symmetric and non-negative definite. Meanwhile, it has other nice properties. For example, given

the tensor sum deposition $\mathcal{H} = \mathcal{H}_0 \oplus \mathcal{H}_1$, we have $R = R_0 + R_1$ where R, R_0 , and R_1 are the RK's of $\mathcal{H}, \mathcal{H}_0$, and \mathcal{H}_1 respectively.

We take the Sobolev space $W_2^m[0, 1]$ as an example (Wang (2011)). Define the reproducing kernel as $R_m(x, z) = \sum_{\nu=0}^m k_\nu(x)k_\nu(z) + (-1)^{m-1}k_{2m}(|x - z|)$ for $x, z \in [0, 1]$, where $k_r(x) = B_r(x)/r!$ are the scaled r th Bernoulli polynomials, with $B_0(x) = 1$, $B'_r(x) = rB_{r-1}(x)$, and $\int_0^1 B_r(x)dx = 0$ for $r = 0, 1, 2, \dots$. The first five scaled Bernoulli polynomials are given below.

$$\begin{aligned} k_0(x) &= 1, \\ k_1(x) &= x - 0.5, \\ k_2(x) &= \frac{1}{2} \left[k_1^2(x) - \frac{1}{12} \right], \\ k_3(x) &= \frac{1}{6} \left[k_1^3(x) - \frac{1}{4}k_1(x) \right], \\ k_4(x) &= \frac{1}{24} \left[k_1^4(x) - \frac{1}{2}k_1^2(x) + \frac{7}{240} \right]. \end{aligned} \tag{1.11}$$

Then for the tensor sum decomposition $W_2^m[0, 1] = \mathcal{H}_0 \oplus \mathcal{H}_1$, where

$$\mathcal{H}_0 = \text{span}\{k_0(x), k_1(x), \dots, k_{m-1}(x)\}, \tag{1.12}$$

$$\mathcal{H}_1 = \left\{ f : \int_0^1 f^{(\nu)}(x)dx = 0, \nu = 0, \dots, m-1, \int_0^1 (f^{(m)}(x))^2 dx < \infty \right\}, \tag{1.13}$$

their corresponding RKs are:

$$R_0(x, z) = \sum_{\nu=0}^{m-1} k_\nu(x)k_\nu(z), \tag{1.14}$$

$$R_1(x, z) = k_m(x)k_m(z) + (-1)^{m-1}k_{2m}(|x - z|), \tag{1.15}$$

under the corresponding inner products:

$$\langle f, g \rangle_0 = \sum_{\nu=0}^{m-1} \int_0^1 f^{(\nu)}(x) dx \int_0^1 g^{(\nu)}(x) dx, \quad (1.16)$$

$$\langle f, g \rangle_1 = \int_0^1 f^{(m)}(x) g^{(m)}(x) dx. \quad (1.17)$$

Thus, the penalty term $\int_0^1 (f^{(m)}(x))^2 dx$ in the objective function (1.8) is the same as $\|P_1 f\|^2$, where P_1 is the projection of $W_2^m[0, 1]$ onto \mathcal{H}_1 . Let $\{\phi_1(x), \dots, \phi_m(x)\}$ be the basis of \mathcal{H}_0 , where $\phi_\nu(x) = k_{\nu-1}(x)$, for $\nu = 1, \dots, m$. Also, define the representers as $\xi_i(x) = R_1(x, x_i)$ for $i = 1, \dots, n$. Then the Kimeldorf–Wahba representer theorem states that the minimizer \hat{f}_λ in (1.8) has the representation

$$\hat{f}_\lambda(x) = \sum_{\nu=1}^m d_\nu \phi_\nu(x) + \sum_{i=1}^n c_i \xi_i(x). \quad (1.18)$$

Based on (1.18), the solution in (1.8) can be written as

$$\hat{\mathbf{f}}_\lambda = \mathbf{T}\mathbf{d} + \mathbf{\Sigma}\mathbf{c}, \quad (1.19)$$

where $\mathbf{T} = \{\phi_\nu(x_i)\}_{i=1, \nu=1}^n$, $\mathbf{\Sigma} = \{R_1(x_i, x_j)\}_{i,j=1}^n$, $\mathbf{d} = (d_1, \dots, d_m)'$, $\mathbf{c} = (c_1, \dots, c_n)'$, $\hat{\mathbf{f}}_\lambda = (\hat{f}_\lambda(x_1), \dots, \hat{f}_\lambda(x_n))'$. Thus, the minimization problem in (1.8) is equivalent to:

$$\operatorname{argmin}_{\mathbf{c}, \mathbf{d}} \frac{1}{n} \|\mathbf{y} - (\mathbf{T}\mathbf{d} + \mathbf{\Sigma}\mathbf{c})\|^2 + \lambda \mathbf{c}' \mathbf{\Sigma} \mathbf{c}. \quad (1.20)$$

The coefficient vectors \mathbf{d} and \mathbf{c} in (1.20) can be determined by solving the following

equations:

$$\begin{aligned}(\boldsymbol{\Sigma} + n\lambda\mathbf{I})\mathbf{c} + \mathbf{T}\mathbf{d} &= \mathbf{y}, \\ \mathbf{T}'\mathbf{c} &= \mathbf{0}.\end{aligned}\tag{1.21}$$

1.2 Smoothing Spline ANOVA

Now consider the multivariate nonparametric regression problem:

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad i = 1, \dots, n,\tag{1.22}$$

where $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$ is the vector of values of d explanatory variables for observation y_i with $x_1 \in \mathcal{X}_1, \dots, x_d \in \mathcal{X}_d$, the independent random errors $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)$ have mean 0 and constant variance σ^2 , and f is defined on the product domain $\mathcal{X} = \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_d$. The domains $\mathcal{X}_1, \dots, \mathcal{X}_d$ are arbitrary sets. Let $\mathcal{H}^{(k)}$ be the RKHS on \mathcal{X}_k with the RK $R^{(k)}$, then there exists a RHKS \mathcal{H} on \mathcal{X} , whose RK equals

$$R(\mathbf{x}, \mathbf{z}) = R^{(1)}(x_1, z_1)R^{(2)}(x_2, z_2) \cdots R^{(d)}(x_d, z_d),\tag{1.23}$$

where $\mathbf{x} = (x_1, \dots, x_d) \in \mathcal{X}$ and $\mathbf{z} = (z_1, \dots, z_d) \in \mathcal{X}$. This \mathcal{H} is called the tensor product RKHS, and is denoted as $\mathcal{H} = \mathcal{H}^{(1)} \otimes \mathcal{H}^{(2)} \otimes \dots \otimes \mathcal{H}^{(d)}$. Then, given the tensor sum decomposition of each marginal space $\mathcal{H}^{(k)} = \mathcal{H}_0^{(k)} \oplus \mathcal{H}_1^{(k)}$ with $\mathcal{H}_0^{(k)} = \text{span}\{1\}$, we have

$$\mathcal{H} = \bigotimes_{k=1}^d \mathcal{H}^{(k)}$$

$$\begin{aligned}
&= \bigotimes_{k=1}^d \left(\mathcal{H}_0^{(k)} \oplus \mathcal{H}_1^{(k)} \right) \\
&= \bigoplus_{B \subseteq \{1, \dots, d\}} \left\{ \left(\bigotimes_{k \in B} \mathcal{H}_1^{(k)} \right) \otimes \left(\bigotimes_{k \in B^c} \mathcal{H}_0^{(k)} \right) \right\} \\
&= \bigoplus_{B \subseteq \{1, \dots, d\}} \mathcal{H}_B \\
&= \mathcal{H}_\emptyset \oplus \sum_{k=1}^d \mathcal{H}_{\{k\}} \oplus \sum_{k < \ell} \mathcal{H}_{\{k, \ell\}} \oplus \dots \oplus \mathcal{H}_{\{1, \dots, d\}}, \tag{1.24}
\end{aligned}$$

where $\mathcal{H}_\emptyset = \bigotimes_{k=1}^d \mathcal{H}_0^{(k)}$. Therefore, for $f \in \mathcal{H}$, it can be decomposed accordingly:

$$f = \mu + \sum_{k=1}^d f_k(x_k) + \sum_{k < \ell} f_{k, \ell}(x_k, x_\ell) + \dots + f_{1, \dots, d}(x_1, \dots, x_d), \tag{1.25}$$

where $\mu \in \mathcal{H}_\emptyset$ is the overall mean, $f_k(x_k) \in \mathcal{H}_{\{k\}}$ is the main effect of x_k , $f_{k, \ell}(x_k, x_\ell) \in \mathcal{H}_{\{k, \ell\}}$ is the two-way interaction between x_k and x_ℓ , and so on. Because the similarity to the classical ANOVA model, (1.25) is called the SS-ANOVA decomposition.

In practice, spaces containing high-order interactions are usually omitted to overcome the curse of dimensionality problem. For example, if we assume $f \in \mathcal{H}_\emptyset \oplus \sum_{k=1}^d \mathcal{H}_{\{k\}}$, i.e., no interactions are considered, then we have the additive model (Hastie and Tibshirani (1990)). An SS-ANOVA model is referred to a model that contains any subset of the components in the SS-ANOVA decomposition. After regrouping terms, the model space of an SS-ANOVA model can be written as

$$\mathcal{Q} = \mathcal{H}^0 \oplus \mathcal{H}^1 \oplus \dots \oplus \mathcal{H}^q, \tag{1.26}$$

where \mathcal{H}^0 contains all functional components that are not penalized, and $\mathcal{H}^1, \dots, \mathcal{H}^q$ are orthogonal RKHS's with RK R^j for $j = 1, \dots, q$. Then for each $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_q)'$, the

estimate \hat{f}_λ in the SS-ANOVA model is

$$\hat{f}_\lambda = \operatorname{argmin}_{f \in \mathcal{Q}} \left(\frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \sum_{j=1}^q \lambda_j \|P_j f\|^2 \right), \quad (1.27)$$

where each \mathcal{H}^j has its own smoothing parameter λ_j , and P_j is the orthogonal projector in \mathcal{Q} onto \mathcal{H}^j . Let $\mathcal{H}_1^* = \mathcal{H}^1 \oplus \cdots \oplus \mathcal{H}^q$, then for $f \in \mathcal{H}_1^*$, we have $f(\mathbf{x}) = \sum_{j=1}^q f_j(\mathbf{x})$ where $f_j \in \mathcal{H}^j$ for $j = 1, \dots, q$. After the reparameterization $\lambda_j \triangleq \lambda/\theta_j$, the minimization problem in (1.27) is equivalent to

$$\hat{f}_{\lambda, \boldsymbol{\theta}} = \operatorname{argmin}_{f \in \mathcal{Q}} \left(\frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda \|P_1^* f\|_*^2 \right), \quad (1.28)$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_q)'$, $P_1^* = \sum_{j=1}^q P_j$ is the orthogonal projector in \mathcal{Q} onto \mathcal{H}_1^* , and \mathcal{H}_1^* is the RKHS with the RK $R_1^* = \sum_{j=1}^q \theta_j R^j$ and inner product $\langle f, g \rangle_* = \sum_{j=1}^q \theta_j^{-1} \langle f_j, g_j \rangle$. Because (1.28) has the same form as (1.8), the same technique can be used to solve $\hat{f}_{\lambda, \boldsymbol{\theta}}$ for fixed λ and $\boldsymbol{\theta}$. The solution $\hat{f}_{\lambda, \boldsymbol{\theta}}$ has a similar representation to (1.18), by the Kimeldorf–Wahba representer theorem:

$$\hat{f}_{\lambda, \boldsymbol{\theta}}(\mathbf{x}) = \sum_{\nu=1}^p d_\nu \phi_\nu(\mathbf{x}) + \sum_{i=1}^n c_i \sum_{j=1}^q \theta_j R^j(\mathbf{x}, \mathbf{x}_i), \quad (1.29)$$

where $\{\phi_1, \dots, \phi_p\}$ is the basis of \mathcal{H}^0 . Again, criteria like GCV can be used to select the tuning parameters λ and $\boldsymbol{\theta}$. See Gu (2013b) and Wang (2011).

1.3 COnponent Selection and Smoothing Operator

For the SS ANOVA model, we can see that the representation in (1.29) contains the representer from each \mathcal{H}^j for $j = 1, \dots, q$. Often it is necessary to perform variable selection and leave out some SS-ANOVA components because the corresponding variables are

not important. This is the motivation behind the COmponent Selection and Smoothing Operator (COSSO) method proposed in Lin and Zhang (2006).

Consider the same tensor product decomposition in (1.24), where each marginal space $\mathcal{H}^{(k)}$ is decomposed into $\{1\} \oplus \mathcal{H}_1^{(k)}$. Then for the model space \mathcal{Q} in (1.26), the null space is $\mathcal{H}^0 = \text{span}\{1\}$. The COSSO procedure estimates $f \in \mathcal{Q}$ using

$$\hat{f}_\tau = \operatorname{argmin}_{f \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \tau^2 \sum_{j=1}^q \|P_j f\|, \quad (1.30)$$

where τ is the tuning parameter. The constant is the only term not penalized. Parametric components such as low order polynomials are also adaptively selected in COSSO. Comparing (1.30) with the SS-ANOVA estimate in (1.27), we can see there are two differences. First, there is only one tuning parameter in (1.30), whereas in (1.27) there are q of them, one for each \mathcal{H}^j . Second, the penalty term in (1.30) is the sum of RKHS norms, but in (1.27) the sum of squared RKHS norms is used instead. The reason that the penalty $\|P_j f\|$ is employed in COSSO is to induce sparsity of the SS-ANOVA components, similar to the L_1 penalty used in the LASSO model (see next section). Indeed, the LASSO is a special case of COSSO when the linear models are considered.

Instead of solving (1.30) directly, the COSSO procedure uses the following equivalent formulation:

$$\hat{f}_\lambda = \operatorname{argmin}_{f \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda_0 \sum_{j=1}^q \theta_j^{-1} \|P_j f\|^2 + \lambda \sum_{j=1}^q \theta_j, \quad (1.31)$$

where $\theta_j \geq 0$ for $j = 1, \dots, q$, λ_0 is a constant, and λ is now the tuning parameter. The advantage of using (1.31) is that it is very similar to the SS-ANOVA formulation in (1.27), except there is an extra penalty on the θ_j 's to induce the sparsity of the components f_j 's. When θ_j 's are fixed, (1.31) is exactly the same minimization problem

as in the SS-ANOVA model. Thus, the COSSO algorithm first initializes $\theta_j = 1$ for $j = 1, \dots, q$, and obtains the SS-ANOVA estimate of f with λ_0 tuned by using either CV or GCV criterion. Denoting this estimate as \hat{f}_{SS} , then it has similar representation as in (1.29):

$$\hat{\mathbf{f}}_{SS} = \hat{d} \mathbf{1}_n + \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \hat{\mathbf{c}}, \quad (1.32)$$

where $\mathbf{1}_n$ is the column vector consisting of n 1's, $\hat{\mathbf{c}} = (\hat{c}_1, \dots, \hat{c}_n)'$, $\boldsymbol{\Sigma}_{\boldsymbol{\theta}} = \sum_{k=1}^q \theta_k \boldsymbol{\Sigma}_k$ with $\boldsymbol{\Sigma}_k = \{R^k(x_i, x_j)\}_{i,j=1}^n$, $\hat{\mathbf{f}}_{SS} = (\hat{f}_{SS}(x_1), \dots, \hat{f}_{SS}(x_n))'$. At the next step, the COSSO algorithm solves the following non-negative garrote problem to obtain the estimates of $\boldsymbol{\theta} = (\theta_1, \dots, \theta_q)'$ with fixed \hat{d} and $\hat{\mathbf{c}}$.

$$\min_{\boldsymbol{\theta}} (\mathbf{z} - \mathbf{G}\boldsymbol{\theta})'(\mathbf{z} - \mathbf{G}\boldsymbol{\theta}) \quad \text{subject to } \theta_j \geq 0, j = 1, \dots, q; \sum_{j=1}^q \theta_j \leq M, \quad (1.33)$$

where $\mathbf{z} = \mathbf{y} - \frac{1}{2}n\lambda_0\hat{\mathbf{c}} - \hat{d}\mathbf{1}_n$, $\mathbf{G} = (\boldsymbol{\Sigma}_1\hat{\mathbf{c}}, \dots, \boldsymbol{\Sigma}_q\hat{\mathbf{c}})$ is an $n \times q$ matrix, and the tuning parameter M is selected by CV or GCV criterion. Denote the estimate of $\boldsymbol{\theta}$ as $\hat{\boldsymbol{\theta}}$. Then at the final step, the COSSO algorithm solves (1.31) again with fixed $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}$, to obtain the final estimate $\hat{\mathbf{f}}_{\hat{\boldsymbol{\theta}}}$. Lin and Zhang (2006) implemented the COSSO using $\mathcal{H}^{(k)} = W_2^2[0, 1]$ for $k = 1, \dots, d$.

1.4 LASSO and Adaptive LASSO

1.4.1 Least Absolute Shrinkage and Selection Operator (LASSO)

Consider the same regression problem as in (1.22) with $f(\mathbf{x}_i) = \sum_{j=1}^d x_{ij}\beta_j$ and $\mathcal{X} = \mathbb{R}^d$. Also assume x_{ij} are standardized so that $\sum_{i=1}^n x_{ij}/n = 0$ and $\sum_{i=1}^n x_{ij}^2 = 1$. Let $\boldsymbol{\beta} = (\beta_1, \dots, \beta_d)'$ be the coefficients vector. Then the LASSO estimate of $\boldsymbol{\beta}$ is the

solution to the following L_1 -penalized regression problem:

$$\hat{\boldsymbol{\beta}}_{\lambda}^{\text{lasso}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left(\frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=1}^d x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^d |\beta_j| \right), \quad (1.34)$$

where λ is the tuning parameter that controls the shrinkage of all the coefficients. As λ increases, LASSO continuously shrinks the coefficients towards zero, and some coefficients are shrunk to exact zero when λ is sufficiently large. The solution in (1.34) has closed form only for some special cases. For example, when there is a single predictor, i.e., $d = 1$, the LASSO solution is a soft-thresholded version of the least square estimate $\hat{\beta}^{\text{ols}}$:

$$\begin{aligned} \hat{\beta}_{\lambda}^{\text{lasso}} &= S(\hat{\beta}^{\text{ols}}, \lambda) \triangleq \operatorname{sign}(\hat{\beta}^{\text{ols}})(|\hat{\beta}^{\text{ols}}| - \lambda)_+ \\ &= \begin{cases} \hat{\beta}^{\text{ols}} - \lambda, & \text{if } \hat{\beta}^{\text{ols}} > 0 \text{ and } \lambda < |\hat{\beta}^{\text{ols}}|, \\ \hat{\beta}^{\text{ols}} + \lambda, & \text{if } \hat{\beta}^{\text{ols}} < 0 \text{ and } \lambda < |\hat{\beta}^{\text{ols}}|, \\ 0, & \text{if } \lambda \geq |\hat{\beta}^{\text{ols}}|, \end{cases} \end{aligned} \quad (1.35)$$

where the least square estimate is $\hat{\beta}^{\text{ols}} = \sum_{i=1}^n x_i y_i$. With multiple predictors, if they are uncorrelated, or equivalently when the design matrix is orthonormal, i.e., $\mathbf{X}'\mathbf{X} = \mathbf{I}$, then once again the LASSO estimate of each β_j is the soft-thresholded version of the least square estimate $\hat{\beta}_j^{\text{ols}}$. For the special case in (1.35), we can see the continuous shrinkage nature of the LASSO estimator, which helps to improve the prediction accuracy due to the bias-variance trade-off. This trade-off is controlled by λ . By allowing small increase in the bias in exchange of substantial decrease in the variance when estimating $\boldsymbol{\beta}$, the prediction error can be reduced substantially.

The pathwise coordinate descent algorithm proposed in Friedman et al. (2007) can be used to efficiently solve (1.34). The idea is to repeatedly apply soft-thresholding with a “partial residual” as a response variable. We can rewrite the right hand side of (1.34)

as:

$$g(\tilde{\beta}) = \frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{k \neq j} x_{ik} \tilde{\beta}_k - x_{ij} \beta_j \right)^2 + \lambda \sum_{k \neq j} |\tilde{\beta}_k| + \lambda |\beta_j|, \quad (1.36)$$

where all the values of β_k for $k \neq j$ are fixed at prior iterations' values $\tilde{\beta}_k(\lambda)$, and we update the estimate of a single coefficient β_j by minimizing (1.36). Based on the solution in (1.35), it is clear that we get

$$\tilde{\beta}_j(\lambda) \leftarrow S \left(\sum_{i=1}^n x_{ij} (y_i - \tilde{y}_i^{(j)}), \lambda \right), \quad (1.37)$$

where $\tilde{y}_i^{(j)} = \sum_{k \neq j} x_{ik} \tilde{\beta}_k$. The update (1.37) is repeated for $j = 1, 2, \dots, d, 1, 2, \dots$ until convergence. The pathwise coordinate descent algorithm returns the LASSO solution for any fixed λ . To tune λ , a criterion such as GCV can be used.

1.4.2 Adaptive LASSO

To evaluate the performance of a variable selection procedure, Fan and Li (2001) proposed the oracle properties. The properties have two parts: one is the consistency of variable selection, i.e., a procedure can correctly identify the true nonzero coefficients; the other is the optimal root-n estimation rate of these true nonzero coefficients. Zou (2006) has shown that the LASSO does not always satisfy the oracle properties, and thus proposed a new procedure called the adaptive LASSO to achieve oracle properties.

The idea of the adaptive LASSO is to penalize the coefficients differently by assigning different weights to them. The adaptive LASSO estimator, $\hat{\beta}^{\text{AL}}$, can be obtained by solving the following L_1 -penalized regression problem:

$$\hat{\beta}_\lambda^{\text{AL}} = \underset{\beta}{\operatorname{argmin}} \left(\frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=1}^d x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^d \hat{w}_j |\beta_j| \right), \quad (1.38)$$

where $\hat{\mathbf{w}} = (\hat{w}_1, \dots, \hat{w}_d)'$ is a data-dependent weights vector defined as $\hat{\mathbf{w}} = 1/|\hat{\boldsymbol{\beta}}|^\gamma$, for a root-n-consistent estimator $\hat{\boldsymbol{\beta}}$ to $\boldsymbol{\beta}$, and a positive γ . Usually we use $\hat{\mathbf{w}} = 1/|\hat{\boldsymbol{\beta}}^{\text{ols}}|^\gamma$.

The algorithm for obtaining adaptive LASSO estimates in (1.38) is simply based on the LASSO algorithm:

The Adaptive LASSO Algorithm

1. Let $x_{ij}^* = x_{ij}/\hat{w}_j$, for $i = 1, \dots, n$, $j = 1, \dots, d$, where $\hat{w}_j = 1/|\hat{\beta}_j^{\text{ols}}|^\gamma$ for fixed γ .
2. Solve the LASSO problem for fixed λ :

$$\hat{\boldsymbol{\beta}}^* = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=1}^d x_{ij}^* \beta_j \right)^2 + \lambda \sum_{j=1}^d |\beta_j|. \quad (1.39)$$

3. The adaptive LASSO estimates for fixed λ and γ are:

$$\hat{\beta}_j^{\text{AL}} = \hat{\beta}_j^*/\hat{w}_j, \quad j = 1, \dots, d. \quad (1.40)$$

The adaptive LASSO has two tuning parameters: λ and γ . To find an optimal pair of (λ, γ) , a method such as cross-validation can be used.

1.5 Hybrid Adaptive Spline

The smoothing spline, smoothing spline ANOVA, and COSSO are regularization approaches. They assume f is in certain model space, and then each solve a regularized least square problem by penalizing the model complexity. In contrast, the regression spline is a basis approximation approach. It starts with representing f in terms of basis functions, such as in (1.2), and then selects the optimal number and types of basis accord-

ing to certain model selection criterion. The hybrid adaptive spline (HAS) is proposed in Luo and Wahba (1997), and it can also be classified as a basis approximation (pursuit) approach. It is called “hybrid” because it is similar to the regression spline, but with smoothing spline representers as basis functions.

When $f \in W_2^m[0, 1]$, the smoothing spline estimator of f has the representation in (1.18), which is a linear combination of the basis in \mathcal{H}_0 and the representers in \mathcal{H}_1 . The HAS idea is to select only k representers adaptively among $\mathcal{R}_1 = \{\xi_1(x), \dots, \xi_n(x)\}$ to approximate f . This is useful especially when the underlying function is spatially inhomogeneous. The tuning parameter k controls the model complexity, and can be selected using the GCV criterion. Luo and Wahba (1997) presented the HAS procedure only for the $m = 2$ case, but the procedure is similar for general m , which is presented as follows.

The HAS procedure

1. *Initialization:* Set $\mathcal{B}_m = \{\phi_1(x), \dots, \phi_m(x)\}$ and let M be an upper bound on the number of basis functions to be selected (including those in \mathcal{B}_m).
2. *Forward Selection:* For $k = m + 1, \dots, M$, select ξ_{j_k} among the representers that have not yet been selected, to maximize the reduction of RSS:

$$\xi_{j_k} = \underset{\xi \in \mathcal{R}_1 \cap \mathcal{B}_{k-1}^c}{\operatorname{argmax}} \left\{ \operatorname{RSS}(\mathcal{B}_{k-1}) - \operatorname{RSS}(\mathcal{B}_{k-1} \cup \{\xi\}) \right\}, \quad (1.41)$$

where \mathcal{B}_{k-1} is the set that contains all the basis functions in \mathcal{B}_m and the $(k - 1 - m)$ representers that have been selected in prior steps, $\operatorname{RSS}(\mathcal{B})$ is the RSS of the ordinary least squares fit of the model with all the basis functions in \mathcal{B} as covariates. Update $\mathcal{B}_k = \mathcal{B}_{k-1} \cup \{\xi_{j_k}\}$.

3. *Elimination:* Select k^* , $m \leq k^* \leq M$, to minimize the GCV criterion:

$$\text{GCV}(k) = \frac{\text{RSS}(\mathcal{B}_k)}{[n - (m + (k - m) \times \text{IDF})]^2}, \quad (1.42)$$

where IDF is the inflated degrees of freedom.

4. *Final Model:* Fit an model with ordinary least squares or ridge regression method using the basis functions contained in \mathcal{B}_{k^*} .

In the calculation of GCV in (1.42), the degrees of freedom of each adaptively selected representers are inflated by a factor of IDF (hence the name inflated degrees of freedom) to account for the adaptive model selection process. Based on simulations, the authors suggested that 1.2 is a good choice for the IDF in HAS, when $m = 2$. A larger IDF may cause fewer basis functions to be selected in the final model. One potential drawback of the HAS procedure is due to this ad hoc choice of IDF. We will delay the discussions of the IDF until Section 1.7. Another drawback of HAS comes from the fact that it uses only one type of basis functions — the smoothing spline representers — in the forward selection. These two drawbacks are the main motivations for the BSML procedure, which we will discuss in Section 1.8. In addition, the search procedure in HAS is greedy which may be trapped to a local optima. We will propose a look-ahead search procedure in Chapter 3 to alleviate this drawback.

1.6 Multivariate Adaptive Regression Splines

Multivariate Adaptive Regression Splines (MARS) is another basis approximation approach that was proposed in Friedman (1991). It assumes that f is a linear combination of tensor products of univariate linear spline basis functions.

Consider the same regression problem as in (1.22) with $\mathcal{X} = \mathbb{R}^d$, and a point $\mathbf{X} \in \mathcal{X}$

is written as $\mathbf{X} = (X_1, \dots, X_d)$. Let \mathcal{T}_1 be the collection of pairs of truncated linear functions:

$$\mathcal{T}_1 = \left\{ (X_j - t)_+, (t - X_j)_+ \right\}_{\substack{t \in \{x_{1j}, x_{2j}, \dots, x_{nj}\} \\ j=1, 2, \dots, d}}, \quad (1.43)$$

where the knot t is one of the observed values of the j th explanatory variable X_j , and X_j ranges over $(-\infty, \infty)$. If all the observed values of the j th explanatory variable are distinct for each $j \in \{1, \dots, d\}$, then there are $2dn$ basis functions in \mathcal{T}_1 . Note that each basis $h(\mathbf{X}) = (X_j - t)_+$ is considered as a function in the entire input space \mathbb{R}^d , although it only depends on a single variable X_j . MARS assumes f has the following representation:

$$f(\mathbf{X}) = \beta_0 + \sum_{m=1}^M \beta_m h_m(\mathbf{X}), \quad (1.44)$$

where $h_m(\mathbf{X})$ is either a function in \mathcal{T}_1 or a product of two or more such functions, and M controls the complexity of f and is user-specified. The advantage of using linear splines in \mathcal{T}_1 as building blocks is that now h_m is nonzero only over the small part of the input space, where all its components are nonzero. This allows the h_m basis functions to operate locally and build up a parsimonious regression surface.

Given the basis functions, the coefficients β_m 's are estimated using ordinary least squares. To construct the set h_m 's, MARS first performs forward selection. Starting with the constant $h_0(\mathbf{X}) = 1$, suppose at the k th step the model set is \mathcal{M}_k which contains all the bases selected in the prior steps, then MARS computes the reduction of RSS by adding to \mathcal{M}_k the products of each basis in \mathcal{M}_k with every pair $\{(X_j - t)_+, (t - X_j)_+\}$ in \mathcal{T}_1 that has not yet been selected, and then selects the pair that minimizes the reduction of RSS. The forward selection stops when adding new terms will make the total number of terms in the model exceed $M + 1$. A backward deletion procedure is then applied, deleting the term one by one whose removal causes the smallest increase in RSS, to produce a sequence of candidate models \hat{f}_λ of each size λ . Denote d as the total number

of terms in the model when the forward selection is completed, where $d \leq M + 1$, then the size $\lambda \in \{d, d - 1, \dots, 2\}$. The final model is \hat{f}_{λ^*} where λ^* is chosen to minimize the GCV criterion. Here are the details of the MARS algorithm.

The MARS procedure

1. *Initialization:* Set $\mathcal{B}_1 = \mathcal{M}_1 = \{1\}$ and $k = 1$. We use \mathcal{B} to denote the set of basis functions selected in the forward selection, and \mathcal{M} to denote the set of model terms h_m 's induced by these basis functions. Let M be the user-specified maximum number of model terms (does not include the constant term) as in (1.44).

2. *Forward Selection:*

while $|\mathcal{M}_k| < M + 1$ **do**

$$\begin{aligned} & \{(X_{j^*} - t^*)_+, (t^* - X_{j^*})_+\} \\ = & \underset{\substack{(X_j - t)_+ \in (\mathcal{T}_1 \cap \mathcal{B}_k^c) \\ (t - X_j)_+ \in (\mathcal{T}_1 \cap \mathcal{B}_k^c)}}{\operatorname{argmax}} \left\{ \operatorname{RSS}(\mathcal{M}_k) - \operatorname{RSS}\left(\mathcal{M}_k \cup \bigcup_{h_\ell \in \mathcal{M}_k} \{(X_j - t)_+ \cdot h_\ell(\mathbf{X}), (t - X_j)_+ \cdot h_\ell(\mathbf{X})\}\right) \right\} \end{aligned}$$

Update:

$$\begin{aligned} \mathcal{B}_{k+1} &= \mathcal{B}_k \cup \{(X_{j^*} - t^*)_+, (t^* - X_{j^*})_+\}, \\ \mathcal{M}_{k+1} &= \mathcal{M}_k \cup \bigcup_{h_\ell \in \mathcal{M}_k} \{(X_{j^*} - t^*)_+ \cdot h_\ell(\mathbf{X}), (t^* - X_{j^*})_+ \cdot h_\ell(\mathbf{X})\}, \\ k &= k + 1. \end{aligned}$$

end while

3. *Backward Elimination:* Let \mathcal{M}_f be the model set returned from the forward

selection. Let $d = |\mathcal{M}_f|$. Set $\mathcal{M}^{(d)} = \mathcal{M}_f$.

for $\lambda = d, d-1, \dots, 2$ **do**

$$\xi^\lambda = \underset{\xi \in \mathcal{M}^{(\lambda)}}{\operatorname{argmin}} \left\{ \operatorname{RSS}(\mathcal{M}^{(\lambda)} \setminus \{\xi\}) - \operatorname{RSS}(\mathcal{M}^{(\lambda)}) \right\}.$$

Set $\mathcal{M}^{(\lambda-1)} = \mathcal{M}^{(\lambda)} \setminus \{\xi^\lambda\}$.

end for

Select λ^* , $1 \leq \lambda^* \leq d$, to minimize the GCV criterion:

$$\operatorname{GCV}(\lambda) = \frac{\sum_{i=1}^n (y_i - \hat{f}_\lambda(\mathbf{x}_i))^2}{[n - (r_\lambda + K_\lambda \times \operatorname{IDF})]^2}, \quad (1.45)$$

where \hat{f}_λ is the ordinary least square fit using the basis functions in $\mathcal{M}^{(\lambda)}$, r_λ is the number of terms in $\mathcal{M}^{(\lambda)}$, K_λ is the number of knots in $\mathcal{M}^{(\lambda)}$, and IDF is the inflated degrees of freedom. For MARS, Friedman suggested to use $\operatorname{IDF} = 3$.

4. *Final Model:* \hat{f}_{λ^*} .

The advantage of MARS is that it can be easily applied to high-dimensional data. Also, it allows high-order interactions among the predictor variables, and is very useful when the true model has a hierarchical structure. However, MARS shares the same disadvantages as HAS. They both use fixed IDF, single class of basis (linear splines in MARS), and have a greedy search procedure.

1.7 Generalized Degrees of Freedom and Covariance Penalty

1.7.1 Ideas Behind the IDF and GDF

We have seen in Sections 1.5 and 1.6 that the IDF plays an important role in constructing the model selection criterion in each of HAS and MARS. The concepts of IDF and generalized degrees of freedom (GDF) are also crucial in understanding other adaptive basis selection procedures such as the BSM procedure (Sklar et al. (2013)) and the look-ahead procedure discussed later. For example, during the forward selection stage of the look-ahead procedure, basis functions are selected by taking into account the IDFs of different libraries. The basis functions are grouped into libraries according to certain properties, such as the form of the basis function or smoothness of this function. Previous literature has suggested that it is advantageous to use different IDFs for basis functions that possess different properties. As mentioned earlier, in the HAS procedure (where the bases are smoothing spline representers), Luo and Wahba (1997) suggested $IDF=1.2$; while for the MARS procedure (which has the truncated linear functions as its basis functions), Friedman (1991) suggested $IDF=3$. The potential problem with fixing the IDF is that even if only one library of basis functions is used, the IDF should change as more and more bases are selected, because the IDF should depend on both the number of bases already in the current model, and how well this model fits the data. Therefore, it is restrictive to fix the IDF for each library throughout the searching process. In addition, the choices of IDF in HAS and MARS are ad hoc. The IDF should be chosen adaptively based on the data. For a particular data set, the degrees of freedom of certain basis functions should be inflated less (making them easier to select) compared to the others, when they possess key properties that lead to better approximation of the true function.

However, the situation may change completely for another data set with a different true function. Therefore, the estimation of the IDFs should be data-driven and adaptive. This is also the idea behind the BSML procedure which estimates IDF values based on an estimate of generalized degrees of freedom.

Generalized degrees of freedom (GDF) was proposed in Ye (1998) for Gaussian data, where it was defined as the the sum of the sensitivities of the fitted values to perturbations in the response. Here we summarize Ye's approach. Let $\mathbf{y} = (y_1, \dots, y_n)'$ be the observed values of a response variable Y , $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))'$ be the values of the true function at the design points $\mathbf{x}_1, \dots, \mathbf{x}_n$. The response vector \mathbf{y} is assumed to follow a Gaussian distribution:

$$\mathbf{y} \sim N(\mathbf{f}, \sigma^2 \mathbf{I}). \quad (1.46)$$

Also, let \hat{f} be the estimate of the true function f , and $\hat{\mathbf{f}} = (\hat{f}(\mathbf{x}_1), \dots, \hat{f}(\mathbf{x}_n))'$. A modeling procedure \mathcal{M} is defined as the mapping from \mathbf{y} to $\hat{\mathbf{f}}$:

$$\mathcal{M} : \mathbf{y} \rightarrow \hat{\mathbf{f}}. \quad (1.47)$$

Then the GDF of \mathcal{M} is defined as

$$\begin{aligned} D(\mathcal{M}) &\triangleq \sum_{i=1}^n \frac{\partial \mathbb{E}_{\mathbf{f}}(\hat{f}(\mathbf{x}_i))}{\partial f(\mathbf{x}_i)} \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov} \left(\hat{f}(\mathbf{x}_i), y_i - f(\mathbf{x}_i) \right) \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov} \left(\hat{f}(\mathbf{x}_i), y_i \right). \end{aligned} \quad (1.48)$$

GDF provides a good measure of the model complexity and reduces to the regular degrees of freedom in Gaussian linear regression with given covariates, but its estimation can be computationally expensive in more complicated models. Ye (1998) proposed the

Monte Carlo algorithm to estimate the GDF when it does not have a known analytical value. First, T replicates of Gaussian noises $\boldsymbol{\delta}_t$ (called the perturbations) are simulated, i.e.,

$$\boldsymbol{\delta}_t = (\delta_{t1}, \dots, \delta_{tn}) \sim N(\mathbf{0}, \tau^2 \mathbf{I}_n), \text{ for } t = 1, \dots, T, \quad (1.49)$$

where $\tau \in [0.5\hat{\sigma}, \hat{\sigma}]$ and $\hat{\sigma}$ is an estimate of σ . Then the modeling procedure \mathcal{M} is applied using each perturbed response $\mathbf{y} + \boldsymbol{\delta}_t$ to obtain the estimate $\hat{f}_{\mathbf{y}+\boldsymbol{\delta}_t}$, for $t = 1, \dots, T$. Finally, for fixed \mathbf{x}_i , for each of $i = 1, \dots, n$, the sensitivity \hat{h}_i is calculated as the least-squares estimate of slope in the following model

$$\hat{f}_{\mathbf{y}+\boldsymbol{\delta}_t}(\mathbf{x}_i) = \hat{\alpha} + \hat{h}_i \delta_{ti}, \quad t = 1, \dots, T, \quad (1.50)$$

and the estimate of GDF is $\hat{D}(\mathcal{M}) = \sum_{i=1}^n \hat{h}_i$. Therefore one needs to fit the model \mathcal{M} T times to compute this estimate of the GDF.

The BSML procedure estimates the GDF of the models at each step while performing the forward selection. Let $\hat{D}(\mathcal{M}_k)$ be the estimated GDF of the model \mathcal{M}_k which contains k basis functions, then the IDF of this model is estimated by

$$\widehat{\text{IDF}}(\mathcal{M}_k) = \frac{\hat{D}(\mathcal{M}_k) - m}{k - m}, \quad (1.51)$$

where the IDFs of the first m bases are fixed at 1.

In the next section, we will explain the relationship between GDF, RSS, and MSE.

1.7.2 The Relationship Between GDF, RSS, and MSE

Consider model (1.22). Let \mathcal{M} be the modeling procedure that leads to an estimate $\hat{\mathbf{f}}$. Define the covariance penalty $C(\mathcal{M})$ as

$$C(\mathcal{M}) \triangleq \mathbb{E}[(\mathbf{y} - \mathbf{f})'(\hat{\mathbf{f}} - \mathbf{f})] = \sum_{i=1}^n \text{Cov}(\hat{f}(\mathbf{x}_i), y_i).$$

Let's take a closer look at the covariance penalty $C(\mathcal{M})$. In the standard linear regression setting, conditional on including all of a fixed set of p covariates (analogous to p basis functions evaluated at n discrete points), the $n \times p$ design matrix \mathbf{X} is fixed. Hence the modeling procedure \mathcal{M} will only include the model estimation part, since no selection of the covariates is needed. Let \mathbf{H} be the corresponding hat matrix, i.e., $\hat{\mathbf{y}} = \hat{\mathbf{f}} = \mathbf{H}\mathbf{y}$. Then $C(\mathcal{M})$ has the following explicit form:

$$\begin{aligned} C(\mathcal{M}) &\triangleq \mathbb{E}[(\mathbf{y} - \mathbf{f})'(\hat{\mathbf{f}} - \mathbf{f})] \\ &= \mathbb{E}\{(\mathbf{y} - \mathbf{f})'\hat{\mathbf{f}}\} - \underbrace{(\mathbb{E}(\mathbf{y}) - \mathbf{f})' \mathbf{f}}_{=0} \\ &= \mathbb{E}\{(\mathbf{y} - \mathbf{f})'\hat{\mathbf{f}}\} \\ &= \mathbb{E}(\boldsymbol{\varepsilon}'\mathbf{H}\mathbf{y}) \\ &= \mathbb{E}(\boldsymbol{\varepsilon}'\mathbf{H}\boldsymbol{\varepsilon}) + \mathbb{E}(\boldsymbol{\varepsilon}'\mathbf{H}\mathbf{f}) \\ &= \text{tr}(\mathbf{H})\sigma^2 \\ &= p\sigma^2, \end{aligned} \tag{1.52}$$

where the second-to-last equality is a result of Theorem 1.5 in Seber and Lee (2003). Hence, the GDF $D(\mathcal{M})$ in this linear regression case can be calculated analytically as:

$$D(\mathcal{M}) = \frac{C(\mathcal{M})}{\sigma^2} = p.$$

Because $D(\mathcal{M})$ is a function of the number of parameters p in the model, it measures the model complexity. The IDF for each basis is $D(\mathcal{M})/p = 1$ here. If, however, the modeling procedure \mathcal{M} also includes the variable (basis) selection part, then the “cost” of selecting a basis function will typically be greater than 1¹. The actual GDF and IDF depend on many factors such as the specific modeling procedure \mathcal{M} , the true function, the nature of variables (basis functions), and which variables (basis functions) have been selected in prior steps. It does not have an explicit form in general.

As shown in Sklar et al. (2013), the expected value of the RSS of an model \mathcal{B} can be decomposed into:

$$\begin{aligned}
 \mathbb{E}(\text{RSS}) &= \mathbb{E}(\|\hat{\mathbf{f}} - \mathbf{y}\|^2) \\
 &= \mathbb{E}(\|\hat{\mathbf{f}} - \mathbf{f} + \mathbf{f} - \mathbf{y}\|^2) \\
 &= \mathbb{E}(\|\hat{\mathbf{f}} - \mathbf{f}\|^2) + \mathbb{E}(\|\mathbf{f} - \mathbf{y}\|^2) - 2\mathbb{E}[(\mathbf{y} - \mathbf{f})'(\hat{\mathbf{f}} - \mathbf{f})] \\
 &= \text{MSE} + \mathbb{E}(\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon}) - 2C(\mathcal{M}) \\
 &= \text{MSE} + n\sigma^2 - 2C(\mathcal{M}).
 \end{aligned}$$

That is,

$$\text{MSE} = \mathbb{E}(\text{RSS}) + 2C(\mathcal{M}) - n\sigma^2. \quad (1.53)$$

On the other hand, there is a well-known bias-variance decomposition for the MSE:

$$\begin{aligned}
 \text{MSE} &= \mathbb{E}(\|\mathbf{f} - \hat{\mathbf{f}}\|^2) \\
 &= \mathbb{E}(\|\mathbf{f} - \mathbb{E}(\hat{\mathbf{f}}) + \mathbb{E}(\hat{\mathbf{f}}) - \hat{\mathbf{f}}\|^2) \\
 &= \|\mathbf{f} - \mathbb{E}(\hat{\mathbf{f}})\|^2 + \mathbb{E}(\|\hat{\mathbf{f}} - \mathbb{E}(\hat{\mathbf{f}})\|^2)
 \end{aligned}$$

¹For the basis functions that are always included in the model, e.g., the first m basis functions in the HAS procedure, the IDF will be 1.

$$= \|\text{Bias}(\hat{\mathbf{f}})\|^2 + \text{tr}(\text{Var}(\hat{\mathbf{f}})), \quad (1.54)$$

where $\text{Bias}(\hat{\mathbf{f}}) = \hat{\mathbf{f}} - \mathbb{E}(\hat{\mathbf{f}})$, and $\text{Var}(\hat{\mathbf{f}})$ is the covariance matrix of $\hat{\mathbf{f}}$. For the second term in (1.54), if $\hat{\mathbf{f}}$ is estimated by a modeling procedure \mathcal{M} that does not include selection, then

$$\begin{aligned} \text{tr}(\text{Var}(\hat{\mathbf{f}})) &= \mathbb{E} \left(\|\hat{\mathbf{f}} - \mathbb{E}(\hat{\mathbf{f}})\|^2 \right) \\ &= \mathbb{E} \left(\|\mathbf{H}\mathbf{y} - \mathbb{E}(\mathbf{H}\mathbf{y})\|^2 \right) \\ &= \mathbb{E} \left(\|\mathbf{H}\mathbf{y} - \mathbf{H}\mathbb{E}(\mathbf{y})\|^2 \right) \\ &= \mathbb{E} \left(\|\mathbf{H}\mathbf{y} - \mathbf{H}\mathbf{f}\|^2 \right) \\ &= \mathbb{E} \left(\|\mathbf{H}\boldsymbol{\varepsilon}\|^2 \right) \\ &= \mathbb{E} (\boldsymbol{\varepsilon}' \mathbf{H}' \mathbf{H} \boldsymbol{\varepsilon}) \\ &= \mathbb{E} (\boldsymbol{\varepsilon}' \mathbf{H} \boldsymbol{\varepsilon}), \quad \text{since } \mathbf{H} \text{ is idempotent,} \\ &= p\sigma^2. \end{aligned} \quad (1.55)$$

We can see that the $C(\mathcal{M})$ in (1.52) is the same as $\text{tr}(\text{Var}(\hat{\mathbf{f}}))$ in (1.55) under the fixed covariates linear regression case. Therefore, (1.53) provides another view of the bias-variance trade-off. If additional basis functions are included in a model, the RSS will be smaller, but on the other hand the model's GDF will increase. An ideal linear regression model should achieve the balance between the two.

1.8 The BSML Procedure

The BSML procedure was proposed in Sklar et al. (2013). It selects basis functions adaptively from multiple libraries, where each library consists of basis functions with

similar forms and properties. Compared to using a single library, the advantage of using multiple libraries is that only relatively few basis functions need to be selected from each library to approximate the target function, particularly if the target function is spatially inhomogeneous and if the basis functions in different libraries capture different inhomogeneous features found in the true function. There are infinitely many choices of the libraries. Libraries may be selected from different families including Fourier, spline, radial, wavelet bases and so on. They may also be selected from different types within a family. We can have B-splines, truncated polynomials and reproducing kernel representers for the spline family. Within each type, we can specify different orders of basis, e.g., linear or cubic for polynomial splines.

Table 1.1 lists some commonly used libraries of basis functions (Sklar et al. (2013)). Here \mathcal{T}_m are truncated polynomials, and the step functions \mathcal{T}_0 is a special case with $m = 0$. $\mathcal{R}_{(1,m)}$ are the polynomial spline representers defined as $\mathcal{R}_{(1,m)} = \{\{R_m(x_i, z_1)\}_{i=1}^n, \dots, \{R_m(x_i, z_q)\}_{i=1}^n\}$ where z_1, \dots, z_q are knots. $R_m(x_i, z_j)$ is the reproducing kernel defined in (1.15). For the family \mathcal{P}_2 , the periodic spline reproducing kernel of order-2 is defined as $R_{per,2}(s, t) = \sum_{v=1}^{\infty} \frac{2}{(2\pi v)^4} \cos 2\pi v(s - t)$ for $s, t \in [0, 1]$.

Table 1.1: Notations for libraries of basis functions that will be used in Section 2.2.

Notation	Family of Basis Functions
\mathcal{U}_0	$\{1\}$, constant functions
\mathcal{U}_m	$\{1, x, x^2, \dots, x^m\}$, polynomial functions of orders $0, 1, 2, \dots, m$
\mathcal{C}_m	$\mathcal{R}_{(1,m)}$, polynomial spline representers
\mathcal{P}_2	$\{R_{per,2}(z_j, x)\}$, periodic smoothing spline representers of order-2
\mathcal{T}_m	$\{(x - z_1)_+^m, (x - z_2)_+^m, \dots, (x - z_d)_+^m\}$, truncated ploynomials
$\mathcal{F}_{f,m}$	$\{\sin(2\pi f k x), \cos(2\pi f k x) : k = 1, \dots, m\}$, Fourier basis functions

The BSML procedure starts with the null library \mathcal{L}_0 , which contains all the basis functions that will be included in the model automatically. Let $m = |\mathcal{L}_0|$ and M be the

pre-specified maximum number of basis functions we want to select (including those in \mathcal{L}_0). Basis functions are selected from L additional libraries $\mathcal{O} = \cup_{l=1}^L \mathcal{L}_l$ one at a time. At each step k , denote the sequentially selected basis functions as ϕ_k for $k = m+1, \dots, M$. Let $\mathcal{B}_k = \{\phi_1, \dots, \phi_k\}$ for $k = m, \dots, M$, where $\mathcal{B}_m = \mathcal{L}_0$. Write “model \mathcal{B}_k ” for “a linear combination of the basis functions in \mathcal{B}_k ”, and also \mathcal{M}_k for the modeling procedure that includes both basis functions selection and estimation steps.

There are two variations of the BSML procedures in Sklar et al. (2013): BSML-C and BSML-S. The basic idea of the BSML-S procedure is as follows: at each step k , it first selects the best basis function from each library \mathcal{L}_l for $l = 1, \dots, L$ to form a candidate set; it then selects the best basis ϕ_k among these candidates according to a criterion such as the DPC in (1.58); finally, it uses either the CIC or GCV criterion defined in (1.60) and (1.61) respectively, to select the best model \mathcal{B}_{k^*} , for $m \leq k^* \leq M$.

The BSML-C procedure is a simplified version of BSML-S. At each step k , it selects the best basis function from the combined non-null libraries \mathcal{O} . The final model \mathcal{B}_{k^*} is chosen to minimize either the CIC or GCV criterion defined in (1.60) and (1.61) respectively, for $m \leq k^* \leq M$. The difference between BSML-C and BSML-S is that BSML-C does not distinguish between the selection costs of basis functions from different non-null libraries in the forward selection, so in general it is outperformed by BSML-S.

The crucial quantity used in computing various selection criteria is the GDF (1.48). In the BSML procedures, the GDF is estimated using the same Monte Carlo method discussed in Section 1.7.1. As mentioned earlier, GDF is a good measure of model complexity, but is computationally costly to estimate when an analytical formula is not known. The BSML-C procedure combines the non-null libraries together, therefore the GDF estimates that need to be computed are $\{\hat{D}(\mathcal{M}_{m+1}), \dots, \hat{D}(\mathcal{M}_M)\}$, where \mathcal{M}_k is the modeling procedure at the k th step of forward selection. Hence, the Monte Carlo algorithm is called $M - m$ times in the BSML-C procedure. For the BSML-S procedure,

however, because it treats different libraries separately, things become more complicated. Two different types of GDFs are computed. One is the basis selection cost within each library, denoted as $\hat{D}(\mathcal{M}_{k,l})$, for $l = 1, \dots, L$, where $\mathcal{M}_{k,l}$ is the modeling procedure that selects the basis function $\psi_{l,j_l}^{(k)}$ in the l th non-null library given the model \mathcal{B}_{k-1} and performs the estimation based on $\mathcal{B}_{k-1} \cup \{\psi_{l,j_l}^{(k)}\}$. The other is the average selection cost for each library, denoted as $\hat{D}(\mathcal{M}_k^{(l)})$, where $\mathcal{M}_k^{(l)}$ is the modeling procedure that selects $k-m$ additional basis functions $\{\varphi_{l,j_l}, j = 1, \dots, k-m\}$ one at a time from library \mathcal{L}_l given the model \mathcal{B}_m and performs the estimation based on $\mathcal{B}_m \cup (\bigcup_{j=1}^{k-m} \{\varphi_{l,j_l}\})$. Hence, with the same upper bound M on the total number of basis functions to be selected, the BSML-S procedure calls the Monte Carlo algorithm $2L \cdot (M - m)$ times. If L is large, the BSML-S procedure becomes very computationally demanding. For this reason, application of the BSML procedure mainly focuses on the univariate case only, with small number of libraries.

Here we give the details of the BSML-S procedure. At step k , where $k \in \{m + 1, \dots, M\}$, the procedure first selects a basis function $\psi_{l,j_l}^{(k)}$ for each $l \in \{1, \dots, L\}$ among the remaining basis functions that have not yet been selected in library \mathcal{L}_l , to maximize the reduction in the RSS:

$$\psi_{l,j_l}^{(k)} = \underset{\psi \in \mathcal{L}_l \cap \mathcal{B}_{k-1}^c}{\operatorname{argmax}} \{ \operatorname{RSS}(\mathcal{B}_{k-1}) - \operatorname{RSS}(\mathcal{B}_{k-1} \cup \{\psi\}) \}. \quad (1.56)$$

Then among these candidate basis functions $\Psi_k = \{\psi_{1,j_1}^{(k)}, \dots, \psi_{L,j_L}^{(k)}\}$, the BSML-S procedure selects the one $\phi_k \in \Psi_k$ to minimize the doubly penalized criterion, i.e.,

$$\phi_k = \underset{\psi_{l,j_l}^{(k)} \in \Psi_k}{\operatorname{argmin}} \operatorname{DPC}(\psi_{l,j_l}^{(k)}), \quad (1.57)$$

and

$$\text{DPC}(\psi_{l,jl}^{(k)}) = \text{RSS}\left(\mathcal{B}_{k-1} \cup \{\psi_{l,jl}^{(k)}\}\right) + c_1 \sigma^2 \hat{D}(\mathcal{M}_{k,l}) + c_2 \sigma^2 \hat{A}_k(\mathcal{L}_l), \quad (1.58)$$

where $\hat{D}(\mathcal{M}_{k,l})$ is the selection cost for $\psi_{l,jl}^{(k)}$ in the l th library at step k , $\hat{A}_k(\mathcal{L}_l)$ is the estimated IDF for the l th library up to step k , and c_1, c_2 are constants. The BSML-S procedure computes $\hat{A}_k(\mathcal{L}_l)$ using

$$\hat{A}_k(\mathcal{L}_l) = \frac{\hat{D}(\mathcal{M}_k^{(l)}) - m}{k - m}, \quad (1.59)$$

where $\hat{D}(\mathcal{M}_k^{(l)})$ is the average selection cost for the l th library up to step k , as introduced previously. Sklar et al. (2013) suggested $c_1 = \ln 2$ and $c_2 = 2$. See the supplement materials of Sklar et al. (2013) for the estimation of σ^2 .

In the elimination step, the covariance inflation criterion (CIC) or the generalized cross-validation (GCV) criterion can be used to select the final model. As shown in (1.53), $\text{MSE}(\mathcal{B}_k) = \mathbb{E}(\text{RSS}(\mathcal{B}_k)) + 2C(\mathcal{M}_k) - n\sigma^2$. Hence, the covariance inflation criterion (CIC) can be defined as follows:

$$\begin{aligned} \text{CIC}(k) &\triangleq \frac{1}{n} \text{RSS}(\mathcal{B}_k) + \frac{2}{n} C(\mathcal{M}_k) \\ &= \frac{1}{n} \text{RSS}(\mathcal{B}_k) + \frac{2}{n} \sigma^2 D(\mathcal{M}_k), \end{aligned} \quad (1.60)$$

which is an unbiased estimate of $\text{MSE}(\mathcal{B}_k) + \sigma^2$. In practice, σ^2 and $D(\mathcal{M}_k)$ are replaced by their estimates. The GCV criterion is defined as

$$\text{GCV}(k) = \frac{\sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i))^2}{(n - D(\mathcal{M}_k))^2} = \frac{\text{RSS}(\mathcal{B}_k)}{(n - D(\mathcal{M}_k))^2}. \quad (1.61)$$

To estimate the GDF $D(\mathcal{M}_k)$ in each of (1.60) and (1.61), the following approximation

is used:

$$\hat{D}(\mathcal{M}_k) \approx m + \sum_{h=m+1}^k \hat{D}(\psi_{l_h, j_{l_h}}^{(h)} | \mathcal{B}_{h-1}) \triangleq m + \sum_{h=m+1}^k \max \left\{ \hat{D}(\mathcal{M}_{h, l_h}) - (h-1), 1 \right\}. \quad (1.62)$$

Putting all things together, the algorithm of the BSML-S procedure is as follows:

The BSML-S procedure

1. *Initialization:* Set $\mathcal{B}_m = \mathcal{L}_0$ and let M be an upper bound on the number of basis functions to be selected (including those in \mathcal{L}_0).
2. *Forward selection:* for $k = m + 1, \dots, M$, do
 - (a) *Select within each library:* for $l = 1, \dots, L$, do
 - i Select $\psi_{l, j_l}^{(k)} \in \mathcal{L}_l$ according to (1.56).
 - ii *Estimate GDFs as needed:* compute $\hat{D}(\mathcal{M}_{k, l})$ and $\hat{D}(\psi_{l, j_l}^{(k)} | \mathcal{B}_{k-1})$.
 - (b) *Select between libraries:* select $\phi_k \in \Psi_k$ to minimize (1.58).
 - (c) *Update:* $\mathcal{B}_k = \mathcal{B}_{k-1} \cup \{\phi_k\}$.
3. *Elimination:* choose $k^*, m \leq k^* \leq M$, as the minimizer of the CIC in (1.60) or the GCV criterion in (1.61).

In comparison, here is the BSML-C procedure:

The BSML-C procedure

1. *Initialization:* Set $\mathcal{B}_m = \mathcal{L}_0$ and let M be an upper bound on the number of basis functions to be selected (including those in \mathcal{L}_0).
2. *Forward selection:* for $k = m + 1, \dots, M$, do

- (a) Select ϕ_k among the basis functions not yet been selected in \mathcal{O} to maximize the reduction in RSS:

$$\phi_k = \underset{\psi \in \mathcal{O} \cap \mathcal{B}_{k-1}^c}{\operatorname{argmax}} \{ \operatorname{RSS}(\mathcal{B}_{k-1}) - \operatorname{RSS}(\mathcal{B}_{k-1} \cup \{\psi\}) \}. \quad (1.63)$$

- (b) Update $\mathcal{B}_k = \mathcal{B}_{k-1} \cup \{\phi_k\}$.

- (c) Estimate $\hat{D}(\mathcal{M}_k)$ and $\hat{C}(\mathcal{M}_k)$.

3. *Elimination:* choose $k^*, m \leq k^* \leq M$, as the minimizer of the CIC in (1.60) or the GCV criterion in (1.61).

Therefore, we can see that in BSML-C all non-null libraries are combined together, so there is no need to select bases and estimate the GDFs separately for each library. This will reduce computation burden, but can sometimes be a disadvantage, especially when the libraries used are very different in nature and need to be treated separately. The performance comparison of BSML-S and BSML-C based on simulations will be presented later in Section 2.1.

Compared to other methods such as HAS and MARS, the better performance of the BSML procedure in terms of the estimation of the true function f , when the true function is spatially inhomogeneous, comes from using the GDF to estimate the costs of selecting additional basis functions, and including multiple libraries so that better basis functions can be selected. However, it still has two drawbacks. First, the search of candidate basis functions at each forward selection step (as given in (1.56) and (1.63)) is greedy. The procedure only allows one candidate basis to be selected after each step. Second, the computation burden of estimating the GDFs is heavy, especially for the BSML-S procedure. In Chapters 3, we develop a new procedure that is less greedy in the forward selection, and at the same time more computationally feasible while still treating bases

from different libraries separately.

Chapter 2

Basis Selection from Multiple Libraries Using Adaptive LASSO

2.1 Adaptive LASSO Basis Selection

As we mentioned in Section 1.4.2, the adaptive LASSO method can be used for variable selection because the L_1 penalty shrinks small coefficients to exact zero. Unlike the regular LASSO method, adaptive LASSO allows different penalties for the coefficients by assigning different weights to their L_1 norms. Consider the same nonparametric regression problem as in (1.22) where $f(\mathbf{x}_i)$ is approximated by a linear combination of basis functions, i.e.,

$$f(\mathbf{x}_i) = \sum_{j=1}^{k^*} \beta_j \xi_j(\mathbf{x}_i), \quad \text{for } \xi_j \in \bigcup_{l=0}^L \mathcal{L}_l, \quad (2.1)$$

where $\{\xi_1, \dots, \xi_m\}$ are the null bases from \mathcal{L}_0 , and $\{\xi_{m+1}, \dots, \xi_{k^*}\}$ are selected from $\{\mathcal{L}_1, \dots, \mathcal{L}_L\}$. Then we may apply the adaptive LASSO method in this basis selection setting, because each basis function needs to be penalized differently.

For the simplicity of discussion, assume that there are 3 libraries: \mathcal{L}_0 , \mathcal{L}_1 , and \mathcal{L}_2 ,

where \mathcal{L}_0 is the null library and all the basis functions in it are automatically included in every model considered. Here is how we apply the adaptive LASSO method to select basis functions. First, a total of M basis functions are preselected in $\mathcal{O} = \cup_{l=1}^2 \mathcal{L}_l$ to maximize the reduction of RSS at each step:

$$\psi^{(k)} = \underset{\psi \in \mathcal{O} \cap \mathcal{B}_{k-1}^c}{\operatorname{argmax}} \{ \operatorname{RSS}(\mathcal{B}_{k-1}) - \operatorname{RSS}(\mathcal{B}_{k-1} \cup \{\psi\}) \}, \quad (2.2)$$

for $k = m+1, \dots, M$, where $\mathcal{B}_m = \mathcal{L}_0$, and $\mathcal{B}_k = \mathcal{B}_{k-1} \cup \{\psi^{(k)}\}$. Let $\mathcal{E} = \{\psi^{(m+1)}, \dots, \psi^{(M)}\}$ be the set of preselected non-null basis functions. Also, let d_k be the reduction of RSS after $\psi^{(k)}$ is added to the model \mathcal{B}_{k-1} , i.e.,

$$d_k = \operatorname{RSS}(\mathcal{B}_{k-1}) - \operatorname{RSS}(\mathcal{B}_{k-1} \cup \{\psi^{(k)}\}), \quad k = m+1, \dots, M. \quad (2.3)$$

Then consider the following objective function:

$$\sum_{i=1}^n \left(y_i - \sum_{k=1}^M \beta_k \psi^{(k)}(\mathbf{x}_i) \right)^2 + \sum_{\{j : \psi^{(j)} \in \mathcal{E} \cap \mathcal{L}_1\}} \lambda_1 d_j^{-\gamma} |\beta_j| + \sum_{\{\ell : \psi^{(\ell)} \in \mathcal{E} \cap \mathcal{L}_2\}} \lambda_2 d_\ell^{-\gamma} |\beta_\ell|, \quad (2.4)$$

where $\{\psi^{(1)}, \dots, \psi^{(m)}\} = \mathcal{L}_0$, and $\lambda_1, \lambda_2, \gamma > 0$. Based on (2.4), we can see that the bases in \mathcal{L}_0 are not penalized, and the penalties of other preselected bases are inversely related to the reduction of RSS caused by adding each basis. Compare to the objective function in (1.38), here we use the weights $\hat{w}_k = d_k^{-\gamma}$, instead of $\hat{w}_k = |\hat{\beta}_k^{\text{ols}}|^{-\gamma}$ as proposed in Zou (2006), because the former works better in our simulations. Note that, in contrary to the linear models, the goal here is to estimate the nonparametric function f rather than the parameters β_k 's. We also tried the SEA-LASSO proposed in Qian and Yang (2013), where the weights in the adaptive LASSO are adjusted by the standard error s_k of $\hat{\beta}_k^{\text{ols}}$, i.e., $\hat{w}_k = (s_k / |\hat{\beta}_k^{\text{ols}}|)^{-\gamma}$, since the basis functions from the same libraries can be highly

correlated. However, the performance is still not as good as that using the reduction of RSS in (2.3) and (2.4). The logic behind assigning weights proportional to the reduction of RSS is analogous to selecting candidate basis functions according to the reduction of RSS. Bases that are important should cause a big decrease in RSS when they are added to the model, so they should be penalized less. On the other hand, the coefficients of bases that make little contribution in reducing the overall RSS will be shrunk to zero. We also penalize non-null libraries differently by assigning tuning parameters to each of them, e.g., λ_1 and λ_2 in (2.4). This controls the selection of basis functions from different libraries when they achieve similar reductions of RSS. The parameters λ_1 and λ_2 are tuned by using the criterion in (2.5).

The upper bound on the total number of candidate models is 2^{M-m} when we consider all the combinations of non-null preselected bases. Let $(\lambda_1, \lambda_2, \gamma)$ be the vector of tuning parameters. Let $\mathcal{M} = \{\mathcal{M}_1, \dots, \mathcal{M}_q\}$ be the set of all candidate models from minimizing (2.4) for $\lambda_1, \lambda_2, \gamma \in \mathbb{R}_+$. Then the final model is chosen to minimize the Bayesian Information Criterion (Schwarz (1978)):

$$\text{BIC}(\mathcal{M}_j) = n \log \text{RSS}(\mathcal{M}_j) + \log(n)|\mathcal{M}_j|, \quad (2.5)$$

where $\text{RSS}(\mathcal{M}_j)$ is the RSS of the ordinary least squares fit of the model using the basis functions in \mathcal{M}_j , and $|\mathcal{M}_j|$ is the number of basis functions contained in \mathcal{M}_j . We used the BIC for model selection as suggested in Qian and Yang (2013). In practice, we only consider a finite grid of points for the three tuning parameters. The BIC is then minimized over all the combinations of tuning parameter values.

The whole procedure is summarized below.

The Adaptive LASSO Basis Selection (ALBS) Procedure

1. *Initialization:* Set $\mathcal{B}_m = \mathcal{L}_0$ and let M be an upper bound on the number of basis functions to be preselected (including those in \mathcal{L}_0).
2. *Preselection:* for $k = m + 1, \dots, M$, do
 - (a) Select $\psi^{(k)}$ according to (2.2). Set $\mathcal{B}_k = \mathcal{B}_{k-1} \cup \{\psi^{(k)}\}$.
 - (b) Compute the reduction of RSS d_k defined in (2.3).
3. *Elimination:* For each combination of tuning parameter values $(\lambda_1, \lambda_2, \gamma)$ in a discrete grid, obtain the candidate model by minimizing the adaptive LASSO objective function in (2.4), and add it to the candidate model set \mathcal{M} . Choose $\mathcal{M}_{j^*} \in \mathcal{M}$ to minimize the BIC in (2.5).

2.2 Simulations to Compare ALBS, BSML, and HAS

We use similar simulation studies as in Sklar et al. (2013) to compare the performances of the ALBS and the BSML procedures. We use the same six univariate test functions and libraries of basis functions, which are given in Table 2.1 and 2.2. We use the same notations as in Table 1.1 to represent the libraries. The test functions are plotted in Figure 2.1.

We can see that all of the test functions are spatially inhomogeneous. It is therefore difficult to estimate them using only one class of basis functions. In particular, the Sine-Jumps, Heavisine, and Blocks-Curves are discontinuous at certain locations. The design points used for all simulations are the grid points $\{x_i = i/n : i = 1, \dots, n\}$, where the sample size is $n = 512$. The knots used in the libraries $\mathcal{P}_2, \mathcal{C}_2, \mathcal{T}_0, \mathcal{T}_2$ are grid points

Table 2.1: Test functions used in simulations.

Function Name	True Function $f(x)$	$SD(f)$
Sine-Jumps	$\sin(2\pi x) - 1_{(0.5,1]}(x) + 1_{(0.25,1]}(x)$	1.003
Heavisine	$2.2[4\sin(4\pi x) - \text{sign}(x - 0.3) - \text{sign}(0.72 - x)]$	6.294
Blocks-Curves	$1_{[0.2,0.4)}(x) + \exp[(x + 0.5)^2]1_{[0.4,0.7)}(x) - x^{10}1_{[0.7,1]}(x)$	1.461
LW6	$\sin[2(4x - 2)] + 2\exp[-256(x - 0.5)^2]$	0.839
LW7	$(4x - 2) + 2\exp[-256(x - 0.5)^2]$	1.265
Poly-Sine	$\sin(16\pi x) - 8(x - 0.5)^2 + 8(x - 0.5)^3 1_{(0.5,1]}(x)$	0.840

Table 2.2: Libraries of basis functions used in simulations. Notations of basis libraries can be found in Table 1.1.

Function Name	Basis Libraries		
	\mathcal{L}_0	\mathcal{L}_1	\mathcal{L}_2
Sine-Jumps	\mathcal{U}_0	\mathcal{P}_2	\mathcal{T}_0
Heavisine	\mathcal{U}_0	\mathcal{P}_2	\mathcal{T}_0
Blocks-Curves	\mathcal{U}_1	\mathcal{C}_2	\mathcal{T}_0
LW6	\mathcal{U}_1	\mathcal{C}_2	\mathcal{T}_2
LW7	\mathcal{U}_1	\mathcal{C}_2	\mathcal{T}_2
Poly-Sine	\mathcal{U}_1	$\mathcal{F}_{8,25}$	\mathcal{T}_2

$\{j/n, j = 4, \dots, n - 3\}$. The response variable is generated according to:

$$y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n, \quad \varepsilon_i, \dots, \varepsilon_n \stackrel{iid}{\sim} N(0, \sigma^2), \quad (2.6)$$

where σ is chosen such that the signal to noise ratio (SNR) defined as $SD(f)/\sigma$ is fixed at 4 for all six test examples, and $SD(f)$ is defined as

$$SD(f) = \sqrt{\frac{1}{n-1} \sum_{i=1}^n \left[f(x_i) - \frac{1}{n} \sum_{i=1}^n f(x_i) \right]^2}. \quad (2.7)$$

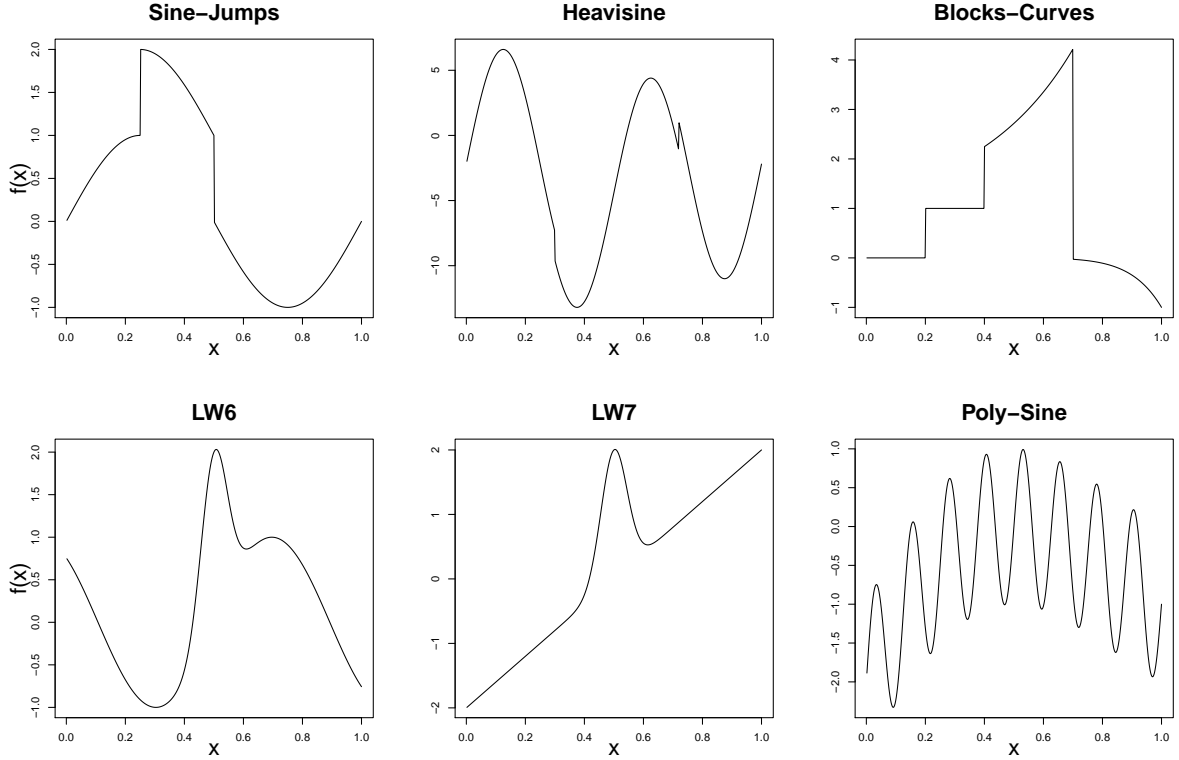


Figure 2.1: Six test functions used in simulations.

We used 100 test sets to compare the performances of four procedures: BSML-S, BSML-C, HAS, and ALBS. Both the HAS and the BSML-C procedures use the combined library $\mathcal{L}_1 \cup \mathcal{L}_2$ in the forward selection, the difference is that BSML-C estimates the GDF by the perturbation method at each step, while HAS fixes the IDF of each selected basis function at 1.2. Although ALBS also carries out the preselection with the combined library, at the following elimination stage it treats the bases from \mathcal{L}_1 and \mathcal{L}_2 differently by imposing different penalties on them. To make a fair comparison between HAS and ALBS, we used the BIC criterion for the model elimination in the HAS procedure. The performance of all procedures is measured by the mean squared error (MSE):

$$\text{MSE}(\hat{f}) = \frac{1}{n} \sum_{i=1}^n \left[\hat{f}(x_i) - f(x_i) \right]^2. \quad (2.8)$$

We used the R package `glmnet` (Friedman et al. (2010)) to obtain the adaptive LASSO estimators in the ALBS procedure. We specified the values each tuning parameter can take as follows:

$$\gamma \in \{e^i : i = -5, -4.5, -4, \dots, 2, 2.5, 3\}, \quad (2.9)$$

$$\frac{\lambda_2}{\lambda_1} \in \{0.1, 0.2, \dots, 0.8, 0.9, 1, 1.5, 2, \dots, 4, 4.5, 5\}, \quad (2.10)$$

and λ_1 values are grid points decided automatically inside the R function `glmnet`, with the default `nlambda = 100`. In the default setting of `glmnet`, the tuning parameter sequence `lambda` is linear on the log scale from `lambda.max` down to `lambda.min`, where `lambda.max` is the smallest value for `lambda` such that all the coefficients are zero, and by default `lambda.min = 0.0001 \times lambda.max`. See the `glmnet` vignette (Hastie and Qian (2014)) for more details. Therefore, the λ_1 values determined inside the `glmnet` function are different for the six test functions, because for each example the `lambda.max` is different.

For each combination of γ , λ_1 , and λ_2 , function `glmnet` returns a candidate model, which is then added into the candidate model set \mathcal{M} . For the other three procedures HAS, BSML-C, and BSML-S, the R package `bsml` (Wu et al. (2012)) was used. For all procedures, the maximal number of bases M was fixed at 30. Figure 2.2 shows the boxplots of the MSEs of 100 test sets using each procedure.

We can see from Figure 2.2 that ALBS had similar performance as HAS for the Blocks-Curves, Heavisine, LW7, and Poly-Sine examples. Although ALBS and HAS preselects exactly the same M basis functions, the set of candidate models considered in ALBS is usually larger than that in HAS. This is because ALBS also penalizes bases differently based on the libraries they are from, so the candidate models returned from the adaptive LASSO algorithm include additional models besides $\mathcal{B}_m, \dots, \mathcal{B}_M$. However, this

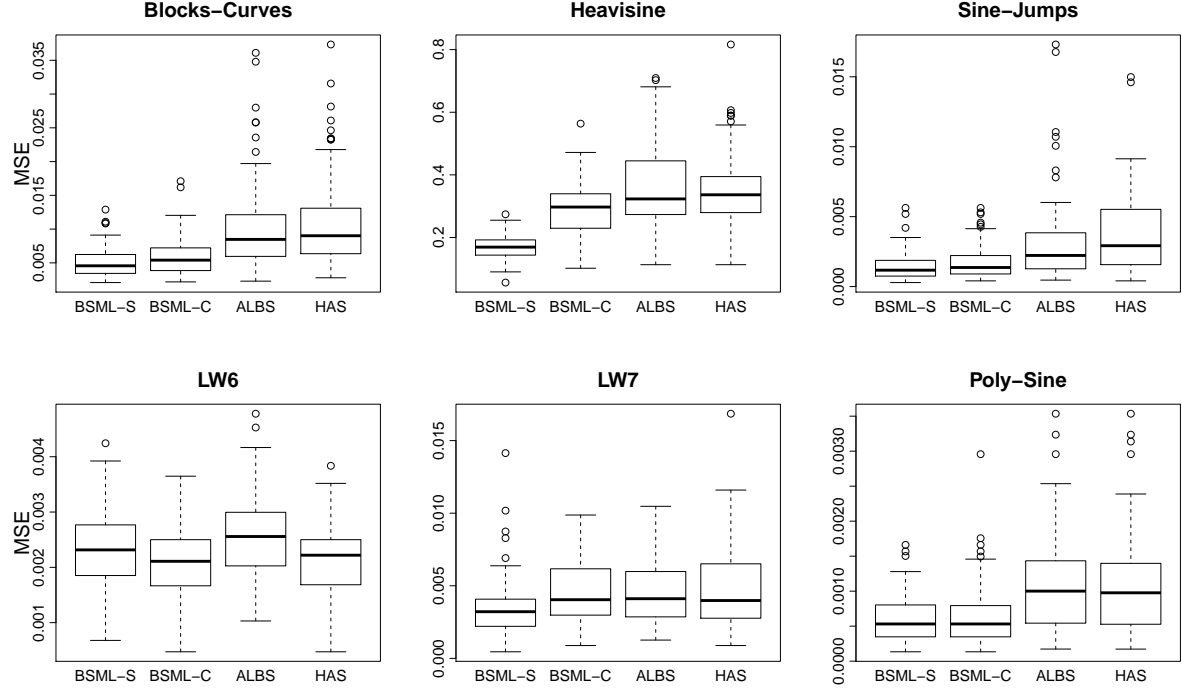


Figure 2.2: Boxplots of the MSEs of 100 test sets using BSML-S, BSML-C, ALBS, and HAS.

did not help the overall performance of the ALBS procedure according to our simulations. ALBS performed worse than the BSML-C procedure (in most cases), especially for the Block-Curves and Heavisine examples. This makes sense because ALBS does not use IDF in the model selection, so in general it tends to select more basis functions than necessary. For example, for the Block-Curves and Heavisine examples, on average BSML-C selected 7 and 12 bases respectively, while ALBS selected 9 and 16 bases respectively. Not surprisingly, the BSML-S procedure had the best overall performance, because it treats the bases from \mathcal{L}_1 and \mathcal{L}_2 differently, both in the forward selection and in the model elimination with adaptively estimated GDFs.

2.3 Adaptive LASSO Basis Selection with Estimated IDF

As we found in the last section, the ALBS procedure did not perform as well as the BSML procedure. One of the reasons is that ALBS does not use IDF in the model selection, so that it can overfit the data by selecting too many basis functions. To overcome this drawback, we propose a slightly different version of the ALBS procedure by incorporating the adaptively estimated IDF for each library.

Now in the preselection, the same perturbation method mentioned in Section 1.7.1 is used to estimate the IDF of each preselected basis. Let $\hat{D}(\mathcal{M}_k)$ be the estimate of GDF for model \mathcal{B}_k . We fixed the IDF of the null bases to be 1, so that $\hat{D}(\mathcal{M}_k) = k$ for $k = 1, \dots, m$. Let $\Psi_l = \mathcal{E} \cap \mathcal{L}_l$ for $l = 1, 2$. Then the IDF of library \mathcal{L}_l is estimated using:

$$\widehat{\text{IDF}}(\mathcal{L}_l) = \max \left(\frac{1}{|\Psi_l|} \sum_{\{j : \psi^{(j)} \in \Psi_l\}} \left(\hat{D}(\mathcal{M}_j) - \hat{D}(\mathcal{M}_{j-1}) \right), 1 \right). \quad (2.11)$$

In the model elimination stage, it uses the following modified BIC to select the final model:

$$\text{BIC}(\mathcal{M}_j) = n \log \text{RSS}(\mathcal{M}_j) + \log(n) \left(m + |\mathcal{M}_j \cap \mathcal{L}_1| \widehat{\text{IDF}}(\mathcal{L}_1) + |\mathcal{M}_j \cap \mathcal{L}_2| \widehat{\text{IDF}}(\mathcal{L}_2) \right). \quad (2.12)$$

The new procedure is referred as the ALBS-2, which is summarized below.

The ALBS-2 Procedure

1. *Initialization:* Set $\mathcal{B}_m = \mathcal{L}_0$ and let M be an upper bound on the number of basis functions to be preselected (including those in \mathcal{L}_0). Set $\hat{D}(\mathcal{M}_k) = k$ for $k = 1, \dots, m$.

2. *Preselection:* for $k = m + 1, \dots, M$, do
 - (a) Select $\psi^{(k)}$ according to (2.2).
 - (b) Compute the reduction of RSS d_k defined in (2.3).
 - (c) Estimate $\hat{D}(\mathcal{M}_k)$ using the perturbation method.
3. *Elimination:* For each combination of tuning parameter values $(\lambda_1, \lambda_2, \gamma)$ in a discrete grid, obtain the candidate model by minimizing the adaptive LASSO objective function in (2.4), add it to the candidate model set \mathcal{M} . Estimate the IDF of each non-null library using (2.11). Choose $\mathcal{M}_{j^*} \in \mathcal{M}$ to minimize the BIC in (2.12).

2.4 Simulations to Compare ALBS-2, BSML, and HAS

We did a similar simulation study as in Section 2.2 to see if ALBS-2 performed better than ALBS. Same 100 test sets were used. The number of perturbations T in the Monte Carlo algorithm was fixed at 50, and the perturbation standard deviation was set at $\tau = 0.5\hat{\sigma}$, where $\hat{\sigma}$ was the Rice estimator of σ (Rice (1984)):

$$\hat{\sigma}_{Rice} = \sqrt{\frac{1}{2(n-1)} \sum_{i=2}^n (y_i - y_{i-1})^2}. \quad (2.13)$$

The boxplots of the MSEs are shown in Figure 2.3. With the adaptively estimated IDFs for each library, ALBS-2 performed better than ALBS in the Blocks-Curves, Sine-Jumps, and Poly-Sine examples. Recall that for the Blocks-Curves and Sine-Jumps examples,

\mathcal{L}_1 and \mathcal{L}_2 contain spline representers and step functions respectively. Thus, for these two examples a bigger IDF for \mathcal{L}_2 relatively to \mathcal{L}_1 helps ALBS-2 to select fewer step functions when choosing the final model. In the Blocks-Curves example, the average IDFs of \mathcal{L}_1 and \mathcal{L}_2 were 2.54 and 3.96 respectively. In the Sine-Jumps example, they were 1.88 and 3.98 respectively. Based on the 100 test sets, on average ALBS selected 4.55 bases from \mathcal{L}_2 for the Sine-Jumps example, and 6.73 bases from \mathcal{L}_2 for the Blocks-Curves example. When ALBS-2 was used, however, these two numbers reduced to 2.03 and 3.08 respectively. The BSML-S procedure was still the best in all examples, but for the Blocks-Curves, Sine-Jumps, and Poly-Sine examples the ALBS-2 procedure performed almost as well. The biggest difference performance-wise between ALBS-2 and BSML-S is with the Heavisine example. The main reason is that for this example, both ALBS-2 and BSML-C combine \mathcal{L}_1 (periodic spline representers) and \mathcal{L}_2 (step functions) in the preselection process, while the BSML-S compares the cost of selecting bases from \mathcal{L}_1 and from \mathcal{L}_2 at every step to decide which library to choose from. For the Heavisine example, when BSML-C, ALBS, and ALBS-2 were used, not enough bases from \mathcal{L}_1 were preselected, and hence fewer were chosen in the final model. This can be seen from Figure 2.4.

The average CPU time of fitting one test set for each procedure and example used is shown in Table 2.3. Overall, the ALBS-2 procedure does not have clear advantages in terms of both performance and speed when compared to the BSML-C procedure. The reason is that they share the same forward selection process, and only differ in how the elimination is carried out.

Although the ALBS and ALBS-2 procedures we proposed so far are only for the special case when there are two non-null libraries, they can be easily extended to general case of L non-null libraries. However, the potential problem with our procedures is that as the number of libraries grows, so does the number of tuning parameters. It

can be easily seen that the number of combinations of gridded parameter values grows exponentially in the number of libraries. Thus, the computation for the ALBS and ALBS-2 procedures is prohibitive when there are too many libraries, for example in the multivariate nonparametric regression case. Another drawback of the ALBS-2 procedure is the way it estimates the IDFs. It is computationally expensive to have to estimate the IDFs at each preselection step, with repeated perturbations. The BSML procedures share the same drawback, especially BSML-S because it needs to estimate the selection cost for each library as well.

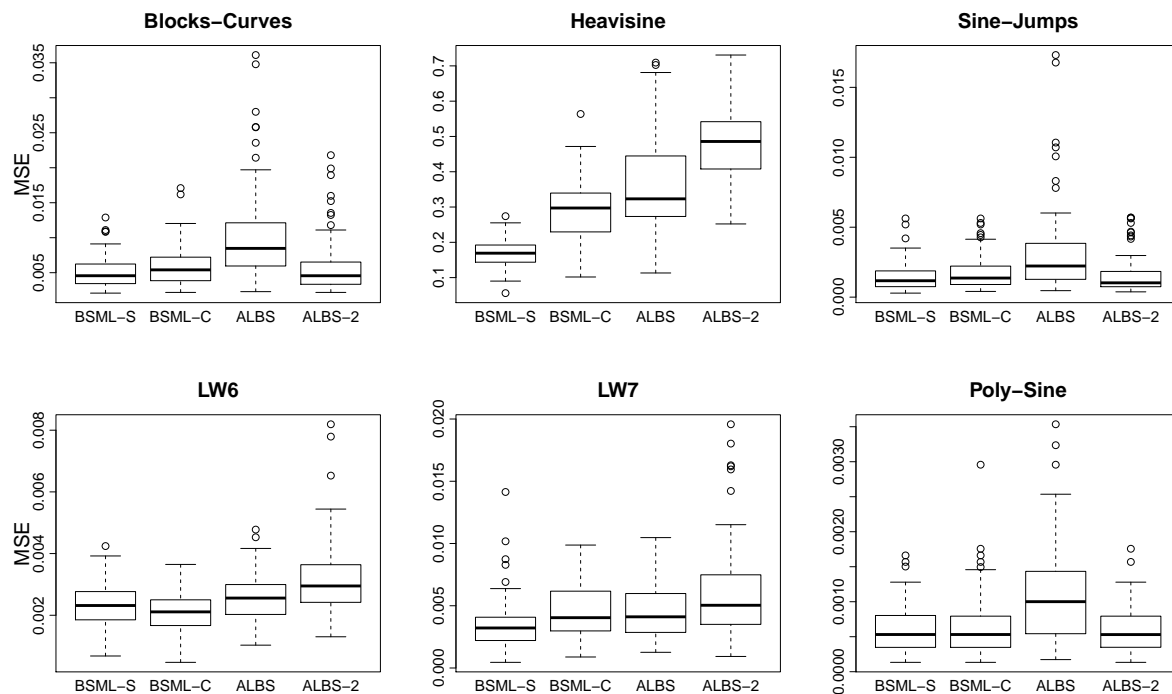


Figure 2.3: Boxplots of the MSEs of 100 test sets using BSML-S, BSML-C, ALBS, and ALBS-2.

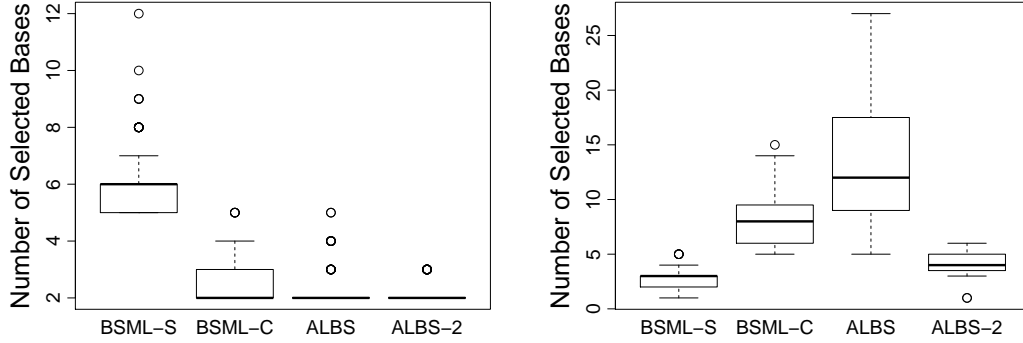


Figure 2.4: Number of selected bases in \mathcal{L}_1 (left) and \mathcal{L}_2 (right) for the Heavisine example with 100 test sets using BSML-S, BSML-C, ALBS, and ALBS-2.

Table 2.3: Average CPU time (in seconds) based on simulations in Sections 2.2 and 2.4.

	Blocks-Curves	Heavisine	LW6	LW7	Poly-Sine	Sine-Jumps
HAS	3.01	2.92	2.94	2.92	0.77	2.93
ALBS	4.13	7.47	7.63	7.78	4.17	4.35
ALBS-2	7.08	9.62	9.31	9.29	5.59	6.70
BSML-C	8.93	9.29	9.26	9.23	4.30	9.22
BSML-S	51.85	51.33	51.70	51.37	28.81	51.20

Based on the simulations results, there are a few things we have learned so far. First of all, each library needs to be assigned with a unique IDF that depends on the data. Also, in the forward selection process, bases from different libraries need to be considered separately. These findings are the main motivations for the look-ahead procedure we will discuss in the next chapter. Because it is costly to estimate the IDFs, the look-ahead procedure instead treats the IDFs as tuning parameters. Meanwhile, the look-ahead procedure evaluates the candidate bases from each individual library in forward selection, and is less greedy in the searching process.

Chapter 3

Look-Ahead Procedure

3.1 Problems With Greedy Search

As mentioned in Section 1.7.1, one of the disadvantages of the BSML procedure is that the GDFs need to be estimated at each step in the forward selection, which can be computationally expensive. Another disadvantage of the BSML procedure comes from its greediness in adding basis functions. At each forward selection step it performs a greedy search for a single basis function that maximizes the reduction in the residual sum of squares, i.e., it only looks one step ahead. This greedy (short-sighted) nature may lead to the selection of a basis function that is not optimal in the long run.

We propose the look-ahead strategy to help alleviate the greediness in the searching for basis functions. The idea of look-ahead is that when choosing the next moves, instead of only evaluating the outcome obtained in a single step, we look further steps ahead and inspect the performance of some extra potential moves. This approach has been widely studied in the context of tree-search and constructive heuristics. See Frost and Dechter (1995), Bertsekas et al. (1997), and Voss et al. (2005). The look-ahead idea has also been applied in sequential Monte Carlo method (e.g., Zhang and Liu (2002)), and in variable

selection, e.g., Zhang et al. (2007), where multiple look-ahead steps are implemented in the best subset selection.

Our look-ahead procedure is less greedy compared to the BSML procedure because it considers more candidate basis functions when performing forward selection. Assuming that there are L non-null libraries, then each time the procedure updates a current model, it considers not only the best single basis function from each library, but also a second basis function among all the bases that have not yet been selected, to form a pair with the best single basis, yielding a total of at most $2L$ candidate bases. Hence, the search space is much bigger than that of the BSML procedure at each update. This look-ahead idea is important particularly in the multiple libraries setting because now it is possible for two basis functions from different libraries to be selected together as a pair. The optimal pair of basis functions may be different if we choose to select them one step at a time, since in that case the selection of the second basis function is conditioning on the previously selected basis function already in the model. Our look-ahead procedure is a compromise between the most greedy one-step ahead procedure (e.g. BSML) and exhaustive search which is computationally prohibitive.

Although one can also modify the forward selection in the BSML procedure to make it look ahead, the increased search space will make the estimation of the GDFs even more computationally expensive. We will develop a computationally feasible procedure that does not require estimating the IDFs.

The look-ahead procedure mainly consists of two parts. The first is the forward selection part and the second is the elimination part. Let's focus on the forward selection part first.

3.2 Forward Selection In the Look-Ahead Procedure with Fixed IDFs

In the forward selection stage of the look-ahead procedure, the procedure takes in some IDF tuning vectors as input and uses them in the calculation of a criterion for building the candidate models. For now, let us focus on the scenario of building a candidate model given a specific IDF tuning vector.

The BSML procedure selects one basis function at each step in the forward selection. In comparison, the look-ahead procedure can select either one or a pair of basis functions to enter the model each time. Suppose that there are in total $L + 1$ libraries, $\{\mathcal{L}_0, \mathcal{L}_1, \dots, \mathcal{L}_L\}$, where \mathcal{L}_0 is the null library with m basis functions that will be included automatically. Let \mathcal{B}_m be the model that includes all the basis functions from the null library, i.e., $\mathcal{B}_m = \mathcal{L}_0$. As the forward selection continues, more and more basis functions are added to \mathcal{B}_m , so we obtain a sequence of nested models. The maximum number of basis functions a model can contain is set to be M ¹. Denote $|S|$ as the cardinality of a set S . Denote \mathcal{B}_- as the model that will be updated next during the iterations, and define index k as $k = |\mathcal{B}_-| + 1$. If there is at least one new basis function being added to \mathcal{B}_- after the update, then the number of basis functions in \mathcal{B}_- will either increase by one if one basis is selected, or increase by two if a pair of basis functions is selected. Correspondingly $|\mathcal{B}_-|$ will increase by either one or two at the next iteration.

We start with $\mathcal{B}_- = \mathcal{L}_0$ and $k = m + 1$. As long as $k < M$, the look-ahead procedure performs the forward selection as follows. Suppose there are L_- libraries among $\{\mathcal{L}_1, \dots, \mathcal{L}_L\}$ that contain at least one unselected basis function, denote them as $\mathcal{L}_{(1)}, \dots, \mathcal{L}_{(L_-)}$, and $\mathcal{O}_- = \cup_{l=1}^{L_-} \mathcal{L}_{(l)}$. The procedure first selects an ordered pair of basis

¹This upper bound is the user-specified variable `MAXBAS` in the code.

functions

$$\psi_l^- = \operatorname{argmax}_{\psi \in \mathcal{L}_{(l)} \cap \mathcal{B}_-^c} \{ \text{RSS}(\mathcal{B}_-) - \text{RSS}(\mathcal{B}_- \cup \{\psi\}) \}, \quad (3.1)$$

and

$$\psi_{\cdot|l}^- = \operatorname{argmax}_{\psi \in \mathcal{O}_- \cap [\mathcal{B}_- \cup \{\psi_l^-\}]^c} \left\{ \text{RSS}(\mathcal{B}_- \cup \{\psi_l^-\}) - \text{RSS}(\mathcal{B}_- \cup \{\psi_l^-, \psi\}) \right\}, \quad (3.2)$$

for each $l \in \{1, 2, \dots, L_-\}$ to form candidates. Let Ψ^- be the collection of all candidate singleton and pairs

$$\Psi^- = \left\{ \{\psi_1^-\}, \dots, \{\psi_{L_-}^-\}, \{\psi_1^-, \psi_{\cdot|1}^-\}, \dots, \{\psi_{L_-}^-, \psi_{\cdot|L_-}^-\} \right\}, \quad (3.3)$$

where the first L_- candidates in (3.3) are single basis functions from (3.1) and the rest are ordered pairs of basis functions. Here some pairs in Ψ^- could be the same except for order, and these redundant pairs are removed from Ψ^- . Then $|\Psi^-|$ satisfies $L_- + 1 \leq |\Psi^-| \leq 2L_-$. Also, note that for the second basis function in these pairs, i.e., $\psi_{\cdot|l}^-$, it is selected given ψ_l^- in the model, and it can be any non-null basis function that has not yet been selected besides ψ_l^- . This selection scheme is different from the one that selects the second basis functions from each of the L_- libraries, i.e., the scheme that finds in total $\frac{1}{2}L_-(L_- + 1)$ pairs to maximize the reduction of RSS,

$$\{\psi_\ell^-, \psi_\kappa^-\} = \left\{ \psi_a, \psi_b : \operatorname{argmax}_{\substack{\psi_a \in \mathcal{L}_{(\ell)} \cap \mathcal{B}_-^c \\ \psi_b \in \mathcal{L}_{(\kappa)} \cap \mathcal{B}_-^c}} \left\{ \text{RSS}(\mathcal{B}_-) - \text{RSS}(\mathcal{B}_- \cup \{\psi_a, \psi_b\}) \right\} \right\}, \quad (3.4)$$

for $1 \leq \ell, \kappa \leq L_-$. Then because $|\mathcal{B}_-| = k - 1$, $|\mathcal{O}_- \cap \{\mathcal{B}_- \cup \{\psi_l^-\}\}^c| = |\mathcal{O}_-| - k + m$, so for the scheme in (3.2), it searches through $|\mathcal{O}_-| - k + m$ basis functions L_- times, or $(|\mathcal{O}_-| - k + m)L_-$ in total. In comparison, for the scheme in (3.4), it searches through

$\binom{|\mathcal{O}_-| - k + m + 1}{2} = \frac{1}{2}(|\mathcal{O}_-| - k + m + 1)(|\mathcal{O}_-| - k + m)$ basis functions. Usually, we have $L_- \ll \frac{1}{2}(|\mathcal{O}_-| - k + m + 1)$, so that the first scheme requires much less computational time. Of course, the trade-off of using (3.2) is that the pairs of basis functions selected may not be the optimal compared to those selected using (3.4). However, based on our simulation experience, the improvement in terms of MSE is too small to justify the extra computations required by the all pairwise search scheme.

After the look-ahead procedure obtains the candidate set Ψ^- in (3.3), we then select $\Psi_{j^*}^- \in \Psi^-$ to maximize the reduction of the BIC. Alternatively, other selection criteria such as the AIC (Akaike (1974)) could also be used. Let $\mathcal{B}_j^- = \mathcal{B}_- \cup \Psi_j^-$ be the j th candidate model for $j = 1, 2, \dots, |\Psi^-|$, where Ψ_j^- is the j th element of Ψ^- :

$$\Psi_j^- = \begin{cases} \{\psi_j^-\}, & \text{if } 1 \leq j \leq L_-, \\ \{\psi_{j-L_-}^-, \psi_{\cdot|j-L_-}^-\}, & \text{if } L_- + 1 \leq j \leq |\Psi^-|. \end{cases} \quad (3.5)$$

The BIC is defined as:

$$\text{BIC}(\mathcal{B}_j^-) = n \log \left(\text{RSS}(\mathcal{B}_j^-) \right) + \log(n) \text{DF}(\mathcal{B}_j^-), \quad (3.6)$$

where $\text{DF}(\mathcal{B}_j^-)$ is the degrees of freedom for model \mathcal{B}_j^- defined as follows:

$$\begin{aligned} \text{DF}(\mathcal{B}_j^-) &= \text{DF}(\mathcal{B}_-) + \text{IDF}(\mathcal{L}(\psi_j^-)) I(1 \leq j \leq L_-) \\ &\quad + \left[\text{IDF}(\mathcal{L}(\psi_{j-L_-}^-)) + \text{IDF}(\mathcal{L}(\psi_{\cdot|j-L_-}^-)) \right] I(L_- + 1 \leq j \leq |\Psi^-|), \end{aligned} \quad (3.7)$$

where $\text{DF}(\mathcal{B}_-)$ is the degrees of freedom of model \mathcal{B}_- , $I(\cdot)$ is the indicator function, $\mathcal{L}(\psi_{\cdot|j-L_-}^-)$ represents the library that the basis function $\psi_{\cdot|j-L_-}^-$ is in. $\text{DF}(\mathcal{B}_-)$ can be

computed using

$$\text{DF}(\mathcal{B}_-) = \text{DF}(\mathcal{B}_m) + \sum_{i=m+1}^{k-1} \text{IDF}(\mathcal{L}(\xi_i)) = m + \sum_{i=m+1}^{k-1} \text{IDF}(\mathcal{L}(\xi_i)), \quad (3.8)$$

where $\xi_{m+1}, \dots, \xi_{k-1}$ are the non-null bases contained in \mathcal{B}_- . Note that $\text{DF}(\mathcal{B}_m) = m$ because the IDF of all the basis functions in the null library \mathcal{L}_0 are fixed at 1.

Let $\Psi_{j^*}^-$ be the set in Ψ^- such that the model $\mathcal{B}_- \cup \Psi_{j^*}^-$ has the lowest BIC, i.e., the index j^* is chosen such that

$$j^* = \underset{1 \leq j \leq |\Psi^-|}{\text{argmin}} \text{BIC}(\mathcal{B}_j^-). \quad (3.9)$$

Alternatively, if AIC is used, (3.6) and (3.9) change into

$$\text{AIC}(\mathcal{B}_j^-) = n \log \left(\text{RSS}(\mathcal{B}_j^-) \right) + 2\text{DF}(\mathcal{B}_j^-) \quad (3.10)$$

$$j^* = \underset{1 \leq j \leq |\Psi^-|}{\text{argmin}} \text{AIC}(\mathcal{B}_j^-). \quad (3.11)$$

Note that if $\text{BIC}(\mathcal{B}_{j^*}^-)$ is greater or equal to $\text{BIC}(\mathcal{B}_-)$, then this means adding the basis set $\Psi_{j^*}^-$ to the current model \mathcal{B}_- will likely cause overfitting, so nothing should be added to \mathcal{B}_- instead, and the forward selection for the current model should be stopped. Otherwise the procedure generates either one or two new candidate models depending on how many bases are selected (corresponding to $\Psi_{j^*}^-$). If $j^* \leq L_-$, then the current model is replaced with $\mathcal{B}_- \cup \{\psi_{j^*}^-\}$. If $j^* > L_-$, then the current model is replaced with $\mathcal{B}_- \cup \{\psi_{j^*-L_-}^-, \psi_{\cdot|j^*-L_-}^-\}$. To start the next iteration, \mathcal{B}_- is set to be the updated model, and $k = |\mathcal{B}_-| + 1$. As long as $k < M$, the same forward selection routine in (3.1), (3.2), (3.5), (3.6), (3.7), and (3.9) is repeated with the new \mathcal{B}_- and k . We can see that different from the BSML procedure, the look-ahead procedure is able to stop the

forward selection before the M th step if the BIC starts to increase. This will help to avoid the overfitting problem when building each model in the forward selection stage and to reduce computation. Another difference of the look-ahead procedure is that it can select two basis functions simultaneously when performing the forward selection, as we have mentioned earlier.

If $k = M$, which means there are already $M - 1$ basis functions in the current model \mathcal{B}_- , the look-ahead procedure will select at most one more basis function so that after this step the number of basis functions in the model will not exceed M . The selection criterion at this step is:

$$\psi_l^- = \underset{\psi \in \mathcal{L}_{(l)} \cap \mathcal{B}_-^c}{\operatorname{argmax}} \left\{ \operatorname{RSS}(\mathcal{B}_-) - \operatorname{RSS}(\mathcal{B}_- \cup \{\psi\}) \right\}, \quad \text{for } 1 \leq l \leq L_-, \quad (3.12)$$

$$l^* = \underset{1 \leq l \leq L_-}{\operatorname{argmin}} \operatorname{BIC}(\mathcal{B}_- \cup \{\psi_l^-\}). \quad (3.13)$$

If $\operatorname{BIC}(\mathcal{B}_- \cup \{\psi_{l^*}^-\}) < \operatorname{BIC}(\mathcal{B}_-)$, then the procedure replaces the current model with $\mathcal{B}_- \cup \{\psi_{l^*}^-\}$.

If $k > M$, the forward selection for the current model will stop.

3.3 Forward Selection Via Householder Transformation

To find the basis function in (3.1), (3.2), and (3.12), we need a fast algorithm to compute the reduction of RSS when one or two new basis functions are added to a model. We are going to use the Householder transformation as in BSML and HAS for selection of a single basis function, and we will extend the computational methods to the case of adding a pair of new basis functions.

Seber and Lee (2003) provides a good discussion of how to use the Householder QR decomposition to calculate the regression coefficients and RSS efficiently. We will summarize some results here.

Consider the linear regression model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (3.14)$$

where $\mathbf{y} = (y_1, \dots, y_n)'$ is the $n \times 1$ response vector, \mathbf{X} is an $n \times k$ design matrix of full column rank ², $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k)'$ is the $k \times 1$ coefficients vector, and $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$ is the $n \times 1$ vector of random errors with zero mean and constant variance. The least squares estimate $\hat{\boldsymbol{\beta}}$ can be obtained by solving the normal equation:

$$\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}. \quad (3.15)$$

After $\hat{\boldsymbol{\beta}}$ is obtained, the RSS can be calculated as:

$$\text{RSS} = \mathbf{e}'\mathbf{e}, \quad (3.16)$$

where $\mathbf{e} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ are the residuals. It is usually computationally expensive to compute $(\mathbf{X}'\mathbf{X})^{-1}$ directly when k is large. The idea is to transform \mathbf{X} into an upper triangular matrix, so that $\hat{\boldsymbol{\beta}}$ in (3.15) can be solved by back substitution.

3.3.1 The Householder QR Decomposition

The Householder transformation can be used to transform the columns of \mathbf{X} so that the elements in the lower triangular will be all zeros. We now describe the Householder transformation. Suppose we have an $n \times 1$ vector $\mathbf{x} = (x_1, x_2, \dots, x_n)'$, and we want to

²In our setting, the columns of \mathbf{X} represent basis functions evaluated at design points.

rotate it so that the 2nd through n th elements become zero in the rotated space without changing the Euclidean norm of \mathbf{x} . Hence, in the rotated space $\tilde{\mathbf{x}} = (\|\mathbf{x}\|, 0, \dots, 0)' = \|\mathbf{x}\|\mathbf{e}_1$, where $\tilde{\mathbf{x}}$ denotes the vector after the rotation, $\|\cdot\|$ is the Euclidean norm, and $\mathbf{e}_1 = (1, 0, \dots, 0)'$ is the n -dimensional vector with only the first element nonzero. Define a Householder vector

$$\mathbf{u} = \mathbf{x} - \|\mathbf{x}\|\mathbf{e}_1 = (x_1 - \|\mathbf{x}\|, x_2, \dots, x_n)', \quad (3.17)$$

then $\mathbf{x} - \mathbf{u}$ will be the vector after the rotation. As mentioned in Seber and Lee (2003), when calculating the first element, $x_1 - \|\mathbf{x}\|$, of \mathbf{u} , severe cancellation would occur if x_1 is positive and large compared with the other elements of \mathbf{x} . Thus, we use a slightly different version of the Householder vector:

$$\begin{aligned} \mathbf{u}^* &= \mathbf{x} + \text{sign}(x_1)\|\mathbf{x}\|\mathbf{e}_1 \\ &= \begin{cases} (x_1 + \|\mathbf{x}\|, x_2, \dots, x_n)', & \text{if } x_1 \geq 0, \\ (x_1 - \|\mathbf{x}\|, x_2, \dots, x_n)', & \text{if } x_1 < 0, \end{cases} \end{aligned}$$

where $\text{sign}(x)$ is 1 if $x \geq 0$, and -1 if $x < 0$. Then the transformed vector becomes:

$$\begin{aligned} \tilde{\mathbf{x}} &= \mathbf{x} - \mathbf{u}^* \\ &= -\text{sign}(x_1)\|\mathbf{x}\|\mathbf{e}_1 \\ &= \begin{cases} (-\|\mathbf{x}\|, 0, \dots, 0)', & \text{if } x_1 \geq 0, \\ (\|\mathbf{x}\|, 0, \dots, 0)', & \text{if } x_1 < 0. \end{cases} \end{aligned} \quad (3.18)$$

To make the Householder matrix defined later to have a simpler form, we scale \mathbf{u}^* so

that it has Euclidean norm 1:

$$\begin{aligned} \mathbf{w} &= \frac{\mathbf{u}^*}{\|\mathbf{u}^*\|} \\ &= \frac{\mathbf{x} + \text{sign}(x_1)\|\mathbf{x}\|\mathbf{e}_1}{\sqrt{2\|\mathbf{x}\|(\|\mathbf{x}\| + |x_1|)}} \end{aligned} \quad (3.19)$$

$$= \begin{cases} \left(\sqrt{\frac{\|\mathbf{x}\| + x_1}{2\|\mathbf{x}\|}}, \frac{x_2}{\sqrt{2\|\mathbf{x}\|(\|\mathbf{x}\| + x_1)}}, \dots, \frac{x_n}{\sqrt{2\|\mathbf{x}\|(\|\mathbf{x}\| + x_1)}} \right)', & \text{if } x_1 \geq 0, \\ \left(-\sqrt{\frac{\|\mathbf{x}\| - x_1}{2\|\mathbf{x}\|}}, \frac{x_2}{\sqrt{2\|\mathbf{x}\|(\|\mathbf{x}\| - x_1)}}, \dots, \frac{x_n}{\sqrt{2\|\mathbf{x}\|(\|\mathbf{x}\| - x_1)}} \right)', & \text{if } x_1 < 0. \end{cases} \quad (3.20)$$

Using $\mathbf{x} - \|\mathbf{u}^*\|\mathbf{w}$ to obtain the vector $\tilde{\mathbf{x}}$ is equivalent to define the following Householder transformation matrix \mathbf{H} , and then applying to \mathbf{x} :

$$\mathbf{H} = \mathbf{I} - 2\mathbf{w}\mathbf{w}'. \quad (3.21)$$

It is easy to see that $\mathbf{H}\mathbf{x} = \tilde{\mathbf{x}}$ by using the facts $\|\mathbf{w}\| = 1$, $\tilde{\mathbf{x}} = \mathbf{x} - \|\mathbf{u}^*\|\mathbf{w}$, and $\|\mathbf{x}\| = \|\tilde{\mathbf{x}}\|$:

$$\begin{aligned} 2\mathbf{w}'\mathbf{x} &= \frac{2}{\|\mathbf{u}^*\|}(\mathbf{x}' - \tilde{\mathbf{x}}')\mathbf{x} \\ &= \frac{2}{\|\mathbf{u}^*\|}(\mathbf{x}'\mathbf{x} - \tilde{\mathbf{x}}'\mathbf{x}) \\ &= \frac{1}{\|\mathbf{u}^*\|}(2\|\mathbf{x}\|^2 - 2\tilde{\mathbf{x}}'\mathbf{x}) \\ &= \frac{1}{\|\mathbf{u}^*\|}(\|\mathbf{x}\|^2 - 2\tilde{\mathbf{x}}'\mathbf{x} + \|\tilde{\mathbf{x}}\|^2) \\ &= \frac{1}{\|\mathbf{u}^*\|}\|\mathbf{x} - \tilde{\mathbf{x}}\|^2 \\ &= \frac{1}{\|\mathbf{u}^*\|}\|\|\mathbf{u}^*\|\mathbf{w}\|^2 \\ &= \|\mathbf{u}^*\|\|\mathbf{w}\|^2 \end{aligned}$$

$$= \|\mathbf{u}^*\|.$$

Consequently,

$$\begin{aligned} \mathbf{H}\mathbf{x} &= \mathbf{x} - 2\mathbf{w}\mathbf{w}'\mathbf{x} \\ &= \mathbf{x} - \mathbf{w}\|\mathbf{u}^*\| \\ &= \tilde{\mathbf{x}}. \end{aligned}$$

The Householder matrix \mathbf{H} has some nice properties. It is obvious the \mathbf{H} is symmetric and orthogonal:

$$\begin{aligned} \mathbf{H}'\mathbf{H} &= \mathbf{H}^2 \\ &= (\mathbf{I} - 2\mathbf{w}\mathbf{w}')(\mathbf{I} - 2\mathbf{w}\mathbf{w}') \\ &= \mathbf{I} - 4\mathbf{w}\mathbf{w}' + 4\mathbf{w}(\mathbf{w}'\mathbf{w})\mathbf{w}' \\ &= \mathbf{I} - 4\mathbf{w}\mathbf{w}' + 4\|\mathbf{w}\|^2\mathbf{w}\mathbf{w}' \\ &= \mathbf{I}. \end{aligned}$$

We can utilize the Householder transformation to perform QR factorization. A rectangular matrix $\mathbf{A} \in \mathbb{R}^{n \times k}$ with $n \geq k$ can be factored into a product of an orthogonal matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ and an upper triangular matrix $\mathbf{U} \in \mathbb{R}^{n \times k}$:

$$\mathbf{A} = \mathbf{Q}\mathbf{U}. \tag{3.22}$$

This factorization is referred to as the QR factorization. Note that if \mathbf{A} has full column rank, then by Theorem 5.2.2 in Golub and Loan (2013) the “diagonal” elements of \mathbf{U} are nonzero, i.e., $\mathbf{U}(j, j) \neq 0$ for $j = 1, \dots, k$. We assume \mathbf{A} to have full column rank from here on. To transform \mathbf{A} into \mathbf{U} , the Householder transformation can be applied to “zero-out” certain elements for the columns in \mathbf{A} . Here is how it works. Denote the

columns of the rank- k $n \times k$ matrix \mathbf{A} as: $\mathbf{A}(:, 1), \mathbf{A}(:, 2), \dots, \mathbf{A}(:, k)$. First, choose \mathbf{H}_1 to zero out all but the first element of $\mathbf{A}(:, 1)$. Following (3.19) and (3.21), let

$$\mathbf{H}_1 = \mathbf{I}_n - 2\mathbf{w}_1\mathbf{w}_1', \quad \text{with } \mathbf{w}_1 = \frac{\mathbf{A}(:, 1) + \text{sign}(\mathbf{A}(1, 1))\|\mathbf{A}(:, 1)\|\mathbf{e}_1}{\sqrt{2\|\mathbf{A}(:, 1)\|(\|\mathbf{A}(:, 1)\| + |\mathbf{A}(1, 1)|)}}, \quad (3.23)$$

where $\mathbf{A}(1, 1)$ is the element in the first row and first column of \mathbf{A} . Premultiplying \mathbf{A} by \mathbf{H}_1 gives

$$\mathbf{H}_1\mathbf{A} = \begin{pmatrix} r_{11} & r_{12} & \cdots & r_{1k} \\ 0 & & & \\ \vdots & \mathbf{A}_1 & & \\ 0 & & & \end{pmatrix},$$

with $r_{11} = -\text{sign}(\mathbf{A}(1, 1))\|\mathbf{A}(:, 1)\|$, and $r_{1j} = \mathbf{H}_1(1, :)\mathbf{A}(:, j)$ for $2 \leq j \leq k$. Next consider a matrix \mathbf{H}_2 of the form

$$\mathbf{H}_2 = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & & & \\ \vdots & \mathbf{K}_2 & & \\ 0 & & & \end{pmatrix},$$

where \mathbf{K}_2 is an $(n-1) \times (n-1)$ Householder matrix chosen to zero out all but the first element of the first column of the submatrix \mathbf{A}_1 , i.e.,

$$\mathbf{K}_2 = \mathbf{I}_{n-1} - 2\mathbf{w}_2\mathbf{w}_2', \quad \text{with } \mathbf{w}_2 = \frac{\mathbf{A}_1(:, 1) + \text{sign}(\mathbf{A}_1(1, 1))\|\mathbf{A}_1(:, 1)\|\mathbf{e}_1}{\sqrt{2\|\mathbf{A}_1(:, 1)\|(\|\mathbf{A}_1(:, 1)\| + |\mathbf{A}_1(1, 1)|)}}, \quad (3.24)$$

We get

$$\mathbf{H}_2 \mathbf{H}_1 \mathbf{A} = \begin{pmatrix} r_{11} & r_{12} & r_{13} & \cdots & r_{1k} \\ 0 & r_{22} & r_{23} & \cdots & r_{2k} \\ 0 & 0 & & & \\ \vdots & \vdots & & \mathbf{A}_2 & \\ 0 & 0 & & & \end{pmatrix},$$

with $r_{22} = -\text{sign}(\mathbf{A}_1(1, 1))\|\mathbf{A}_1(:, 1)\|$, and $r_{2j} = \mathbf{K}_2(1, :)\mathbf{A}_1(:, j - 1)$ for $3 \leq j \leq k$. If we continue analogous steps, the last premultiplication matrix is:

$$\mathbf{H}_k = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & 1 & 0 \\ 0 & \cdots & 0 & \mathbf{K}_k \end{pmatrix},$$

where \mathbf{K}_k is an $(n - k + 1) \times (n - k + 1)$ Householder matrix chosen so that the last $n - k$ elements in the k th column of $\mathbf{H}_k \cdots \mathbf{H}_1 \mathbf{A}$ are all zeros. Therefore, we have:

$$\mathbf{H}_k \mathbf{H}_{k-1} \cdots \mathbf{H}_1 \mathbf{A} = \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} r_{11} & r_{12} & r_{13} & \cdots & r_{1k} \\ 0 & r_{22} & r_{23} & \cdots & r_{2k} \\ \vdots & 0 & r_{33} & \cdots & r_{3k} \\ \vdots & \vdots & 0 & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & r_{kk} \\ \hline \vdots & \vdots & \vdots & \vdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

Here \mathbf{R} is an upper-triangular matrix. We define \mathbf{H}_j as

$$\mathbf{H}_j = \begin{cases} \mathbf{H}_1, & \text{if } j = 1, \\ \begin{pmatrix} \mathbf{I}_{j-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_j \end{pmatrix}, & \text{if } 2 \leq j \leq k. \end{cases} \quad (3.25)$$

Note that $\mathbf{H}_j^2 = \mathbf{H}_j' \mathbf{H}_j = \mathbf{I}_n$ for $j = 1, \dots, k$, because of the symmetric and orthogonal properties of the Householder matrices. If we define $\mathbf{Q} = (\mathbf{H}_k \mathbf{H}_{k-1} \cdots \mathbf{H}_1)' = \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_k$, then we have the Householder QR decomposition:

$$\mathbf{A} = \mathbf{Q} \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix}, \quad (3.26)$$

since

$$\begin{aligned} \mathbf{A} &= \mathbf{I}_n \mathbf{A} \\ &= \mathbf{H}_1^2 \mathbf{A} \\ &= \mathbf{H}_1 \mathbf{I}_n \mathbf{H}_1 \mathbf{A} \\ &= \mathbf{H}_1 \mathbf{H}_2 \mathbf{H}_2 \mathbf{H}_1 \mathbf{A} \\ &= \mathbf{H}_1 \mathbf{H}_2 \mathbf{I}_n \mathbf{H}_2 \mathbf{H}_1 \mathbf{A} \\ &= \mathbf{H}_1 \mathbf{H}_2 \mathbf{H}_3 \mathbf{H}_3 \mathbf{H}_2 \mathbf{H}_1 \mathbf{A} \\ &= \cdots \\ &= (\mathbf{H}_1 \cdots \mathbf{H}_k)(\mathbf{H}_k \cdots \mathbf{H}_1 \mathbf{A}) \\ &= \mathbf{Q} \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix}. \end{aligned}$$

Note that \mathbf{Q} is orthogonal, because it is a product of orthogonal matrices \mathbf{H}_j 's. Further-

more, \mathbf{Q} can be written as $\mathbf{Q} = (\mathbf{Q}_k, \mathbf{Q}_{n-k})$, where \mathbf{Q}_k is $n \times k$ and \mathbf{Q}_{n-k} is $n \times (n-k)$, so that

$$\mathbf{A} = (\mathbf{Q}_k, \mathbf{Q}_{n-k}) \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix} = \mathbf{Q}_k \mathbf{R}.$$

For the linear regression setting as in (3.14), let \mathbf{A} be the $n \times k$ design matrix \mathbf{X} with full column rank, i.e., $\mathbf{A} = \mathbf{X}$, then by (3.26) for the augmented matrix (\mathbf{X}, \mathbf{y}) we have:

$$(\mathbf{X}, \mathbf{y}) = (\mathbf{Q}_k, \mathbf{Q}_{n-k}) \begin{pmatrix} \mathbf{R} & \mathbf{r}_1 \\ \mathbf{0} & \mathbf{r}_2 \end{pmatrix}, \quad (3.27)$$

where \mathbf{R} is a $k \times k$ upper triangular matrix, \mathbf{r}_1 is a $k \times 1$ vector, and \mathbf{r}_2 is a $(n-k) \times 1$ vector, with

$$\begin{pmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \end{pmatrix} = \mathbf{Q}' \mathbf{y} = \mathbf{H}_k \mathbf{H}_{k-1} \cdots \mathbf{H}_1 \mathbf{y},$$

i.e., $(\mathbf{r}'_1, \mathbf{r}'_2)'$ is obtained by premultiplying the transformation matrices $\mathbf{H}_1, \dots, \mathbf{H}_k$ to \mathbf{y} . Multiplying out and equating both sides of (3.27) gives

$$\mathbf{X} = \mathbf{Q}_k \mathbf{R}, \quad (3.28)$$

$$\mathbf{y} = \mathbf{Q}_k \mathbf{r}_1 + \mathbf{Q}_{n-k} \mathbf{r}_2. \quad (3.29)$$

Thus, the solution $\hat{\boldsymbol{\beta}}$ to the normal equation (3.15) is

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} \\ &= (\mathbf{R}' \mathbf{Q}'_k \mathbf{Q}_k \mathbf{R})^{-1} \mathbf{R}' \mathbf{Q}_k (\mathbf{Q}_k \mathbf{r}_1 + \mathbf{Q}_{n-k} \mathbf{r}_2) \\ &= (\mathbf{R}' \mathbf{R})^{-1} \mathbf{R}' \mathbf{r}_1 \\ &= \mathbf{R}^{-1} \mathbf{r}_1, \end{aligned} \quad (3.30)$$

where we used the facts that $\mathbf{Q}'_k \mathbf{Q}_k = \mathbf{I}_k$ and $\mathbf{Q}'_k \mathbf{Q}_{n-k} = \mathbf{0}$ because $\mathbf{Q} = (\mathbf{Q}_k, \mathbf{Q}_{n-k})$ is

an orthogonal matrix. Equation (3.30) means $\hat{\beta}$ is the solution of $\mathbf{R}\hat{\beta} = \mathbf{r}_1$. Because \mathbf{R} is upper triangular, $\hat{\beta}$ can be solved by back-substitution fairly easily:

$$\begin{cases} \hat{\beta}_k = \mathbf{r}_1(k)/\mathbf{R}(k, k), \\ \hat{\beta}_j = [\mathbf{r}_1(j) - \sum_{j'=j+1}^k \mathbf{R}(j, j')\hat{\beta}_{j'}] / \mathbf{R}(j, j), \quad \text{for } j = k-1, \dots, 2, 1. \end{cases} \quad (3.31)$$

In addition, the residuals \mathbf{e} and the RSS in (3.16) can be written as:

$$\begin{aligned} \mathbf{e} &= \mathbf{y} - \mathbf{X}\hat{\beta} \\ &= \mathbf{Q}_k \mathbf{r}_1 + \mathbf{Q}_{n-k} \mathbf{r}_2 - \mathbf{Q}_k \mathbf{R} \mathbf{R}^{-1} \mathbf{r}_1 \\ &= \mathbf{Q}_{n-k} \mathbf{r}_2 \\ &= (\mathbf{Q}_k, \mathbf{Q}_{n-k}) \begin{pmatrix} \mathbf{0} \\ \mathbf{r}_2 \end{pmatrix} \\ &= \mathbf{Q} \begin{pmatrix} \mathbf{0} \\ \mathbf{r}_2 \end{pmatrix}, \end{aligned}$$

so that

$$\begin{aligned} \text{RSS} &= \|\mathbf{e}\|^2 \\ &= (\mathbf{0}', \mathbf{r}_2') \mathbf{Q}' \mathbf{Q} (\mathbf{0}', \mathbf{r}_2')' \\ &= (\mathbf{0}', \mathbf{r}_2') (\mathbf{0}', \mathbf{r}_2')' \\ &= \mathbf{r}_2' \mathbf{r}_2. \end{aligned} \quad (3.32)$$

3.3.2 Adding A Single Basis Function Using Householder QR Decomposition

Using the Householder QR decomposition, we can easily compute the reduction of RSS as in (3.1) and select a single candidate basis function from each library. Suppose for the current model \mathcal{B}_- , the basis functions (including the m bases in \mathcal{L}_0) are $\mathcal{B}_- = \{\phi_1, \dots, \phi_{k-1}\}$ where $k > m$. Let $\mathbf{X}_{k-1} = (\mathbf{x}_1, \dots, \mathbf{x}_{k-1})$ be the design matrix corresponding to the basis functions in \mathcal{B}_- , where the column vector \mathbf{x}_j is the basis function ϕ_j evaluated at the design points for each $j \in \{1, \dots, k-1\}$. Let \mathbf{H}_j be the transformation matrix as defined in (3.25), and

$$\mathbf{Q} = \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_{k-1}, \quad (3.33)$$

then similar to (3.27) we have

$$(\mathbf{X}_{k-1}, \mathbf{y}) = \mathbf{Q} \begin{pmatrix} \mathbf{R} & \mathbf{r}_1 \\ 0 & r_2 \\ \mathbf{0} & \mathbf{r}_3 \end{pmatrix}, \quad (3.34)$$

where $(\mathbf{r}'_1, r_2, \mathbf{r}'_3)' = \mathbf{Q}'\mathbf{y}$. Note that here \mathbf{r}_1 and \mathbf{r}_3 are vectors with length $k-1$ and $n-k$ respectively, while r_2 is a scalar. Now we want to add a new basis function ψ to \mathcal{B}_- , which corresponds to adding a new column \mathbf{x} to \mathbf{X}_{k-1} . Let $\mathbf{q} = \mathbf{Q}'\mathbf{x} = (\mathbf{q}'_1, q_2, \mathbf{q}'_3)'$ where \mathbf{q}'_1, q_2 , and \mathbf{q}'_3 have the same dimensions as \mathbf{r}'_1, r_2 , and \mathbf{r}'_3 respectively. Then we have

$$(\mathbf{X}_{k-1}, \mathbf{x}, \mathbf{y}) = \mathbf{Q} \begin{pmatrix} \mathbf{R} & \mathbf{q}_1 & \mathbf{r}_1 \\ 0 & q_2 & r_2 \\ \mathbf{0} & \mathbf{q}_3 & \mathbf{r}_3 \end{pmatrix}. \quad (3.35)$$

To compute the reduction of RSS caused by adding this new basis function, we can apply a Householder transformation to zero out \mathbf{q}_3 . Let \mathbf{K} be the Householder matrix for $\mathbf{c} = (q_2, \mathbf{q}_3')'$, and let $\mathbf{c}^* = \mathbf{K}\mathbf{c} = (q_2^*, \mathbf{0}')'$. Let $\mathbf{d} = (r_2, \mathbf{r}_3')'$ and $\mathbf{d}^* = \mathbf{K}\mathbf{d} = (r_2^*, \mathbf{r}_3^{*'})'$. Then by the result in (3.32), we have

$$\begin{aligned} \text{RSS}(\mathcal{B}_-) &= \mathbf{d}'\mathbf{d} = (\mathbf{K}'\mathbf{d}^*)'(\mathbf{K}'\mathbf{d}^*) = (\mathbf{d}^*)'\mathbf{K}\mathbf{K}'\mathbf{d}^* = (\mathbf{d}^*)'\mathbf{d}^* \\ &= (r_2^*)^2 + (\mathbf{r}_3^{*'})'(\mathbf{r}_3^*) \\ &= (r_2^*)^2 + \text{RSS}(\mathcal{B}_- \cup \{\psi\}). \end{aligned} \quad (3.36)$$

Therefore, the reduction of RSS in (3.1) is:

$$\Delta\text{RSS} = \text{RSS}(\mathcal{B}_-) - \text{RSS}(\mathcal{B}_- \cup \{\psi\}) = (r_2^*)^2. \quad (3.37)$$

However, it is not computationally efficient to use (3.37) directly to compute the reduction of RSS when ψ is added to \mathcal{B}_- . The reason is that in order to find ψ_l^- in (3.1), we need to compute ΔRSS for each of the individual basis functions $\psi \in \mathcal{L}_{(l)} \cap \mathcal{B}_-^c$, but the Householder matrix for each of them is different, so if we use (3.37) directly we need to construct the Householder transformation $|\mathcal{L}_{(l)} \cap \mathcal{B}_-^c|$ number of times, for each $l \in \{1, 2, \dots, L_-\}$.

There is a more convenient way to compute ΔRSS without applying the Householder transformation for ψ explicitly. Denote $\mathbf{K}(:, 1)$ as \mathbf{k} , and $\mathbf{K}(:, 2 : (n - k + 1))$ as \mathbf{K}_r , then we have:

$$(\mathbf{c}, \mathbf{d}) = (\mathbf{k}, \mathbf{K}_r) \begin{pmatrix} q_2^* & r_2^* \\ \mathbf{0} & \mathbf{r}_3^{*'} \end{pmatrix}. \quad (3.38)$$

Multiplying out and equating both sides of (3.38) gives

$$\mathbf{c} = q_2^* \mathbf{k},$$

$$\mathbf{d} = r_2^* \mathbf{k} + \mathbf{K}_r \mathbf{r}_3^*.$$

Therefore,

$$\begin{aligned} \mathbf{c}' \mathbf{d} &= q_2^* \mathbf{k}' (r_2^* \mathbf{k} + \mathbf{K}_r \mathbf{r}_3^*) \\ &= q_2^* r_2^* + q_2^* \mathbf{k}' \mathbf{K}_r \mathbf{r}_3^* \\ &= q_2^* r_2^*, \end{aligned}$$

$$\begin{aligned} \mathbf{c}' \mathbf{c} &= (q_2^* \mathbf{k}') (q_2^* \mathbf{k}) \\ &= (q_2^*)^2, \end{aligned}$$

where we used the fact that the Householder matrix $\mathbf{K} = (\mathbf{k}, \mathbf{K}_r)$ is orthogonal. Hence, according to (3.37) we have

$$\Delta \text{RSS} = (r_2^*)^2 = \frac{(\mathbf{c}' \mathbf{d})^2}{\mathbf{c}' \mathbf{c}} = \frac{(q_2^* r_2^* + \mathbf{q}_3' \mathbf{r}_3)^2}{q_2^2 + \mathbf{q}_3' \mathbf{q}_3}. \quad (3.39)$$

If we use (3.39) to compute ΔRSS , then there is no need to apply any Householder transformation for the basis functions that have not yet been selected, i.e., $\psi \in \mathcal{L}_{(l)} \cap \mathcal{B}_-^c$, in finding ψ_l^- for each $l \in \{1, 2, \dots, L_-\}$. Therefore, for each $l \in \{1, 2, \dots, L_-\}$, to find ψ_l^- in $\mathcal{L}_{(l)} \cap \mathcal{B}_-^c$ that maximizes the reduction of RSS, first the Householder transformation based on the basis functions already in \mathcal{B}_- , namely $\{\phi_1, \dots, \phi_{k-1}\}$, is used to obtain the quantities q_2, \mathbf{q}_3, r_2 , and \mathbf{r}_3 . Then ΔRSS in (3.39) is computed for all unselected basis functions in $\mathcal{L}_{(l)} \cap \mathcal{B}_-^c$, and the one with the largest reduction of RSS is chosen to be candidate basis function ψ_l^- . Suppose ϕ_k is the basis function selected among all the candidates at this iteration. At the next iteration, the Householder transformation for ϕ_k is applied to obtain the transformation matrix \mathbf{H}_k . The design matrix becomes

$\mathbf{X}_k = (\mathbf{X}_{k-1}, \mathbf{x}_k)$, where \mathbf{x}_k is ϕ_k evaluated at the design points. The matrix \mathbf{Q} at this iteration becomes $\mathbf{Q} = \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_k$. ΔRSS in (3.39) can then be computed similarly with the updated quantities q_2, \mathbf{q}_3, r_2 , and \mathbf{r}_3 , and used to select $\{\psi_1^-, \dots, \psi_{L-}^-\}$ at this iteration.

3.3.3 Adding A Pair of Basis Functions Using Householder QR Decomposition

We use the same idea as in the previous section to compute the reduction of RSS when a pair of basis functions are added to the model \mathcal{B}_- . Let $\{\psi_1, \psi_2\}$ be any pair of new basis functions, and let \mathbf{x} and \mathbf{z} be ψ_1 and ψ_2 evaluated at the design points respectively. Using the same definition of \mathbf{Q} in (3.33), then similar to (3.34), we have

$$(\mathbf{X}_{k-1}, \mathbf{y}) = \mathbf{Q} \begin{pmatrix} \mathbf{R} & \mathbf{r}_1 \\ 0 & r_2 \\ 0 & r_3 \\ \mathbf{0} & \mathbf{r}_4 \end{pmatrix}, \quad (3.40)$$

where $(\mathbf{r}'_1, r_2, r_3, \mathbf{r}'_4)' = \mathbf{Q}'\mathbf{y}$, the column vectors \mathbf{r}_1 and \mathbf{r}_4 have length $k-1$ and $n-k-1$ respectively, and both r_2 and r_3 are scalars. Let $\mathbf{p} = \mathbf{Q}'\mathbf{x} = (\mathbf{p}'_1, p_2, p_3, \mathbf{p}'_4)'$ and $\mathbf{q} = \mathbf{Q}'\mathbf{z} = (\mathbf{q}'_1, q_2, q_3, \mathbf{q}'_4)'$ where \mathbf{p}_1 and \mathbf{q}_1 are column vectors with length $k-1$, and \mathbf{p}_4 and \mathbf{q}_4 are column vectors with length $n-k-1$, then similar to (3.35) we have

$$(\mathbf{X}_{k-1}, \mathbf{x}, \mathbf{z}, \mathbf{y}) = \mathbf{Q} \begin{pmatrix} \mathbf{R} & \mathbf{p}_1 & \mathbf{q}_1 & \mathbf{r}_1 \\ 0 & p_2 & q_2 & r_2 \\ 0 & p_3 & q_3 & r_3 \\ \mathbf{0} & \mathbf{p}_4 & \mathbf{q}_4 & \mathbf{r}_4 \end{pmatrix}. \quad (3.41)$$

Denote $\mathbf{b} = (p_2, p_3, \mathbf{p}'_4)'$, $\mathbf{c} = (q_2, q_3, \mathbf{q}'_4)'$, and $\mathbf{d} = (r_2, r_3, \mathbf{r}'_4)'$. Let \mathbf{K} be the Householder matrix for \mathbf{b} , then we have

$$\mathbf{K}(\mathbf{b}, \mathbf{c}, \mathbf{d}) = (\mathbf{K}\mathbf{b}, \mathbf{K}\mathbf{c}, \mathbf{K}\mathbf{d}) = \begin{pmatrix} p_2^* & q_2^* & r_2^* \\ 0 & q_3^* & r_3^* \\ \mathbf{0} & \mathbf{q}_4^* & \mathbf{r}_4^* \end{pmatrix}. \quad (3.42)$$

Denote $\mathbf{c}^* = (q_3^*, (\mathbf{q}_4^*)')'$, and $\mathbf{d}^* = (r_3^*, (\mathbf{r}_4^*)')'$. Let \mathbf{K}^* be the Householder matrix for \mathbf{c}^* , then we have

$$\mathbf{K}^*(\mathbf{c}^*, \mathbf{d}^*) = (\mathbf{K}^*\mathbf{c}^*, \mathbf{K}^*\mathbf{d}^*) = \begin{pmatrix} q_3^{**} & r_3^{**} \\ \mathbf{0} & \mathbf{r}_4^{**} \end{pmatrix}. \quad (3.43)$$

From (3.42), we have

$$\mathbf{K}\mathbf{d} = \begin{pmatrix} r_2^* \\ \mathbf{d}^* \end{pmatrix}, \quad \text{or} \quad \mathbf{d} = \mathbf{K}' \begin{pmatrix} r_2^* \\ \mathbf{d}^* \end{pmatrix}.$$

From (3.43), we have

$$\mathbf{K}^*\mathbf{d}^* = \begin{pmatrix} r_3^{**} \\ \mathbf{r}_4^{**} \end{pmatrix}, \quad \text{or} \quad \mathbf{d}^* = (\mathbf{K}^*)' \begin{pmatrix} r_3^{**} \\ \mathbf{r}_4^{**} \end{pmatrix}.$$

Thus, using the result in (3.32) and also the facts that $\mathbf{K}\mathbf{K}' = \mathbf{I}$ and $\mathbf{K}^*(\mathbf{K}^*)' = \mathbf{I}$, we have

$$\begin{aligned} \text{RSS}(\mathcal{B}_-) &= \mathbf{d}'\mathbf{d} \\ &= \begin{pmatrix} r_2^* & (\mathbf{d}^*)' \end{pmatrix} \mathbf{K}\mathbf{K}' \begin{pmatrix} r_2^* \\ \mathbf{d}^* \end{pmatrix} \\ &= (r_2^*)^2 + (\mathbf{d}^*)'\mathbf{d}^* \\ &= (r_2^*)^2 + \begin{pmatrix} r_3^{**} & (\mathbf{r}_4^{**})' \end{pmatrix} \mathbf{K}^*(\mathbf{K}^*)' \begin{pmatrix} r_3^{**} \\ \mathbf{r}_4^{**} \end{pmatrix} \\ &= (r_2^*)^2 + (r_3^{**})^2 + (\mathbf{r}_4^{**})'\mathbf{r}_4^{**} \end{aligned}$$

$$= (r_2^*)^2 + (r_3^{**})^2 + \text{RSS}(\mathcal{B}_- \cup \{\psi_1, \psi_2\})$$

Therefore, the reduction of RSS by adding $\{\psi_1, \psi_2\}$ to model \mathcal{B}_- is:

$$\Delta \text{RSS} = \text{RSS}(\mathcal{B}_-) - \text{RSS}(\mathcal{B}_- \cup \{\psi_1, \psi_2\}) = (r_2^*)^2 + (r_3^{**})^2. \quad (3.44)$$

Same as before, we want to avoid applying the Householder transformations for ψ_1 and ψ_2 explicitly. Denote $\mathbf{K}(:, 1)$ as \mathbf{k}_1 , $\mathbf{K}(:, 2)$ as \mathbf{k}_2 , and $\mathbf{K}(:, 3 : (n - k + 1))$ as \mathbf{K}_r , then we have

$$(\mathbf{b}, \mathbf{c}, \mathbf{d}) = (\mathbf{k}_1, \mathbf{k}_2, \mathbf{K}_r) \begin{pmatrix} p_2^* & q_2^* & r_2^* \\ 0 & q_3^* & r_3^* \\ \mathbf{0} & \mathbf{q}_4^* & \mathbf{r}_4^* \end{pmatrix}, \quad (3.45)$$

Multiplying out and equating both sides of (3.45) gives

$$\begin{cases} \mathbf{b} = p_2^* \mathbf{k}_1, \end{cases} \quad (3.46)$$

$$\begin{cases} \mathbf{c} = q_2^* \mathbf{k}_1 + q_3^* \mathbf{k}_2 + \mathbf{K}_r \mathbf{q}_4^*, \end{cases} \quad (3.47)$$

$$\begin{cases} \mathbf{d} = r_2^* \mathbf{k}_1 + r_3^* \mathbf{k}_2 + \mathbf{K}_r \mathbf{r}_4^*. \end{cases} \quad (3.48)$$

Therefore,

$$\begin{cases} \mathbf{b}'\mathbf{b} = (p_2^*)^2, \end{cases} \quad (3.49)$$

$$\begin{cases} \mathbf{b}'\mathbf{c} = p_2^* q_2^*, \end{cases} \quad (3.50)$$

$$\begin{cases} \mathbf{b}'\mathbf{d} = p_2^* r_2^*, \end{cases} \quad (3.51)$$

$$\begin{cases} \mathbf{c}'\mathbf{c} = (q_2^*)^2 + (q_3^*)^2 + (\mathbf{q}_4^*)' \mathbf{q}_4^* = (q_2^*)^2 + (\mathbf{c}^*)' \mathbf{c}^*, \end{cases} \quad (3.52)$$

$$\begin{cases} \mathbf{c}'\mathbf{d} = q_2^* r_2^* + q_3^* r_3^* + (\mathbf{q}_4^*)' \mathbf{r}_4^* = q_2^* r_2^* + (\mathbf{c}^*)' \mathbf{d}^*, \end{cases} \quad (3.53)$$

where we have used the facts that $(\mathbf{k}_1)' \mathbf{k}_1 = (\mathbf{k}_2)' \mathbf{k}_2 = 1$, $(\mathbf{K}_r)' \mathbf{K}_r = \mathbf{I}$, $(\mathbf{k}_1)' \mathbf{k}_2 = 0$, and $(\mathbf{k}_1)' \mathbf{K}_r = (\mathbf{k}_2)' \mathbf{K}_r = \mathbf{0}$ because the Householder matrix $\mathbf{K} = (\mathbf{k}_1, \mathbf{k}_2, \mathbf{K}_r)$ is orthogonal.

Similarly, denote $\mathbf{K}^*(:, 1)$ as \mathbf{k}^* , and $\mathbf{K}^*(:, 2 : (n - k))$ as \mathbf{K}_r^* , then we have

$$(\mathbf{c}^*, \mathbf{d}^*) = (\mathbf{k}^*, \mathbf{K}_r^*) \begin{pmatrix} q_3^{**} & r_3^{**} \\ \mathbf{0} & \mathbf{r}_4^{**} \end{pmatrix}. \quad (3.54)$$

Multiplying out and equating both sides of (3.54) gives

$$\begin{cases} \mathbf{c}^* = q_3^{**} \mathbf{k}^*, & (3.55) \end{cases}$$

$$\begin{cases} \mathbf{d}^* = r_3^{**} \mathbf{k}^* + \mathbf{K}_r^* \mathbf{r}_4^{**}. & (3.56) \end{cases}$$

Therefore,

$$\begin{cases} (\mathbf{c}^*)' \mathbf{c}^* = (q_3^{**})^2, & (3.57) \end{cases}$$

$$\begin{cases} (\mathbf{c}^*)' \mathbf{d}^* = q_3^{**} r_3^{**}, & (3.58) \end{cases}$$

where $(\mathbf{k}^*)' \mathbf{k}^* = 1$, $(\mathbf{K}_r^*)' \mathbf{K}_r^* = \mathbf{I}$, and $(\mathbf{k}^*)' \mathbf{K}_r^* = \mathbf{0}$ because of the orthogonality of the Householder matrix $\mathbf{K}^* = (\mathbf{k}^*, \mathbf{K}_r^*)$. By (3.49) and (3.51), we have:

$$(r_2^*)^2 = \frac{(\mathbf{b}' \mathbf{d})^2}{\mathbf{b}' \mathbf{b}}. \quad (3.59)$$

By (3.49) - (3.51),

$$(q_2^*)^2 = \frac{(\mathbf{b}' \mathbf{c})^2}{\mathbf{b}' \mathbf{b}}, \quad q_2^* r_2^* = \frac{(\mathbf{b}' \mathbf{c})(\mathbf{b}' \mathbf{d})}{\mathbf{b}' \mathbf{b}}. \quad (3.60)$$

Thus, using (3.52), (3.53), and (3.60), we have

$$\begin{cases} (\mathbf{c}^*)' \mathbf{c}^* = \mathbf{c}' \mathbf{c} - (q_2^*)^2 = \mathbf{c}' \mathbf{c} - \frac{(\mathbf{b}' \mathbf{c})^2}{\mathbf{b}' \mathbf{b}}, & (3.61) \end{cases}$$

$$\begin{cases} (\mathbf{c}^*)' \mathbf{d}^* = \mathbf{c}' \mathbf{d} - q_2^* r_2^* = \mathbf{c}' \mathbf{d} - \frac{(\mathbf{b}' \mathbf{c})(\mathbf{b}' \mathbf{d})}{\mathbf{b}' \mathbf{b}}. & (3.62) \end{cases}$$

Hence, by (3.57), (3.58), and also the results in (3.61) and (3.62), we have

$$\begin{aligned}
 r_3^{**} &= \frac{((\mathbf{c}^*)' \mathbf{d}^*)^2}{(\mathbf{c}^*)' \mathbf{c}^*} \\
 &= \frac{\left(\mathbf{c}' \mathbf{d} - \frac{(\mathbf{b}' \mathbf{c})(\mathbf{b}' \mathbf{d})}{\mathbf{b}' \mathbf{b}} \right)^2}{\mathbf{c}' \mathbf{c} - \frac{(\mathbf{b}' \mathbf{c})^2}{\mathbf{b}' \mathbf{b}}} \\
 &= \frac{\left[(\mathbf{b}' \mathbf{b})(\mathbf{c}' \mathbf{d}) - (\mathbf{b}' \mathbf{c})(\mathbf{b}' \mathbf{d}) \right]^2}{(\mathbf{b}' \mathbf{b}) \left[(\mathbf{b}' \mathbf{b})(\mathbf{c}' \mathbf{c}) - (\mathbf{b}' \mathbf{c})^2 \right]}. \tag{3.63}
 \end{aligned}$$

Therefore, by (3.44), (3.59), and (3.63), we have

$$\begin{aligned}
 \Delta \text{RSS} &= (r_2^*)^2 + (r_3^{**})^2 \\
 &= \frac{(\mathbf{b}' \mathbf{d})^2}{\mathbf{b}' \mathbf{b}} + \frac{\left[(\mathbf{b}' \mathbf{b})(\mathbf{c}' \mathbf{d}) - (\mathbf{b}' \mathbf{c})(\mathbf{b}' \mathbf{d}) \right]^2}{(\mathbf{b}' \mathbf{b}) \left[(\mathbf{b}' \mathbf{b})(\mathbf{c}' \mathbf{c}) - (\mathbf{b}' \mathbf{c})^2 \right]} \\
 &= \frac{(\mathbf{b}' \mathbf{d})^2(\mathbf{c}' \mathbf{c}) + (\mathbf{b}' \mathbf{b})(\mathbf{c}' \mathbf{d})^2 - 2(\mathbf{c}' \mathbf{d})(\mathbf{b}' \mathbf{c})(\mathbf{b}' \mathbf{d})}{(\mathbf{b}' \mathbf{b})(\mathbf{c}' \mathbf{c}) - (\mathbf{b}' \mathbf{c})^2}, \tag{3.64}
 \end{aligned}$$

where $\mathbf{b} = (p_2, p_3, \mathbf{p}'_4)'$, $\mathbf{c} = (q_2, q_3, \mathbf{q}'_4)'$, and $\mathbf{d} = (r_2, r_3, \mathbf{r}'_4)'$ as defined earlier.

Similar to the situation in Section 3.3.2, to add a pair of basis functions to \mathcal{B}_- , first the Householder transformation based on the basis functions already in \mathcal{B}_- , namely $\{\phi_1, \dots, \phi_{k-1}\}$, is used to obtain the quantities \mathbf{b} , \mathbf{c} , and \mathbf{d} for each pair of basis functions that are considered. Then ΔRSS in (3.64) is computed for each pair of basis functions, and the pair leading to the largest reduction of RSS is selected to be candidate pair. In the look-ahead procedure, a candidate pair has the form $\{\psi_l^-, \psi_{\cdot|l}^-\}$, for $l \in 1, 2, \dots, L_-$, where $\psi_l^- \in \mathcal{L}_{(l)} \cap \mathcal{B}_-^c$. Suppose the pair $\{\phi_k, \phi_{k+1}\}$ is selected among all the candidate basis functions at this iteration. At the next iteration, the Householder transformations for ϕ_k and ϕ_{k+1} are applied to obtain the transformation matrices \mathbf{H}_k and \mathbf{H}_{k+1} respectively.

The design matrix becomes $\mathbf{X}_{k+1} = (\mathbf{X}_{k-1}, \mathbf{x}_k, \mathbf{x}_{k+1})$, where \mathbf{x}_k and \mathbf{x}_{k+1} are ϕ_k and ϕ_{k+1} evaluated at the design points respectively. The matrix \mathbf{Q} at this iteration becomes $\mathbf{Q} = \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_k \mathbf{H}_{k+1}$. ΔRSS in (3.64) can then be computed similarly with the updated quantities \mathbf{b} , \mathbf{c} , and \mathbf{d} , and used to select the candidate pairs at this iteration.

3.4 Forward Selection In the Look-Ahead Procedure with IDFs as Tuning Parameters

We have discussed how the look-ahead procedure performs the forward selection with fixed IDFs in Section 3.2. However, as mentioned in Section 1.7.1, the IDFs should be chosen adaptively based on the data. We now treat the IDFs as tuning parameters and discuss how this will change the forward selection procedure.

Suppose there are $L + 1$ libraries $\{\mathcal{L}_0, \mathcal{L}_1, \dots, \mathcal{L}_L\}$ where \mathcal{L}_0 is the null library and the rest are non-null libraries. The IDF of \mathcal{L}_0 is fixed at 1. The IDFs of the non-null libraries are treated as tuning parameters, which can take any positive values, with $\text{IDF} = 1$ meaning the degrees of freedom is not inflated. Let s_l be the IDF of the l -th library. We consider a user-specified finite number of possible values for s_l and denote the collection of these values as S_l , for $l = 1, \dots, L$. Define the IDF tuning vector as $\mathbf{s} = (s_1, \dots, s_L)$. Then \mathbf{s} is in the IDF space $\mathcal{S} = S_1 \times \cdots \times S_L$, which contains all possible combinations of the user-specified IDF values. Obviously, the size of \mathcal{S} , $\prod_{l=1}^L |S_l|$, grows exponentially in L .

The difference when multiple IDF tuning vectors (several $\mathbf{s} \in \mathcal{S}$) are used is that now when $\text{BIC}(\mathcal{B}_j^-)$ is computed using (3.6), it can attain different minimums depending on the involved IDFs. Hence, there can be multiple choices of $\Psi_{j^*}^-$ according to (3.9), each minimizing $\text{BIC}(\mathcal{B}_j^-)$ given a specific combination of IDFs (i.e., given a specific $\mathbf{s} \in \mathcal{S}$).

As a result, at each step the forward selection process can generate multiple new models $\mathcal{B}_- \cup \Psi_{j^*}^-(\mathbf{s})$ for $\mathbf{s} \in \mathcal{S}$, given the current model \mathcal{B}_- .

During the forward selection process, we need to keep the records of both the basis functions contained in each model and the IDF tuning vectors that lead to the selection of these basis functions. Let \mathcal{D} be the collection of all models generated at the current stage in the forward selection. Here we refer to a unique collection of basis functions as a model. Denote the h th model in \mathcal{D} as \mathcal{A}_h , and let \mathcal{S}_h be the set of IDF tuning vectors that leads to the selection of the basis functions in \mathcal{A}_h . We will call \mathcal{S}_h the IDF set of \mathcal{A}_h from now on. Let \mathcal{V}_- be the set of models in \mathcal{D} that the procedure is updating at the current round³. At this round, the procedure updates one model in \mathcal{V}_- at each iteration. Note that \mathcal{V}_- stays the same at each round of iterations. Let \mathcal{V}_+ be the set of models will be updated at the next round. Before the current round of iterations begins, $\mathcal{V}_- = \mathcal{V}_+$. Then after each model in \mathcal{V}_- is updated, \mathcal{V}_+ is adjusted accordingly. If the procedure decides to stop adding more bases to a set in \mathcal{V}_- , that model is removed from \mathcal{V}_+ . If there are new models generated, they are added to \mathcal{V}_+ and \mathcal{D} . Once all models in \mathcal{V}_- have been updated, the procedure sets $\mathcal{V}_- = \mathcal{V}_+$ for the next round of iterations. We start with a single model $\mathcal{D} = \mathcal{V}_- = \mathcal{V}_+ = \{\mathcal{A}_1\}$, where \mathcal{A}_1 contains the m bases from \mathcal{L}_0 , and we initialize $\mathcal{S}_1 = \mathcal{S}$. Denote the current model being updated by the forward selection process as \mathcal{B}_- , and its IDF set as \mathcal{S}_- . They are initialized with $\mathcal{B}_- = \mathcal{A}_1$ and $\mathcal{S}_- = \mathcal{S}_1$.

Suppose \mathcal{D} , \mathcal{V}_- contain sets of selected basis functions $\mathcal{A}_1, \dots, \mathcal{A}_d$, where d is the total number of models generated so far. At the first iteration, $d = 1$. Suppose now the procedure is updating \mathcal{A}_h in \mathcal{V}_- , where \mathcal{A}_h contains $k - 1$ basis functions. Thus, we have $\mathcal{B}_- = \mathcal{A}_h$, with the IDF set $\mathcal{S}_- = \mathcal{S}_h$. Our procedure then follows the similar forward selection routine in (3.1), (3.2), (3.5), (3.7), and (3.9) for each $\mathbf{s} \in \mathcal{S}_-$. Recall that after

³See the example at the end of this section for more details about each “round” of iterations.

each update, the procedure can add zero, one, or two basis functions to \mathcal{B}_- . For any $\mathbf{s} \in \mathcal{S}_-$, if $\text{BIC}_{\mathbf{s}}(\mathcal{B}_-)$ is smaller or equal to the BICs of all the candidate models, then no basis will be added to \mathcal{B}_- . Define

$$\mathcal{S}_{h,0} = \left\{ \mathbf{s} \in \mathcal{S}_- : \text{BIC}_{\mathbf{s}}(\mathcal{B}_-) \leq \min_{1 \leq j \leq |\Psi^-|} \text{BIC}_{\mathbf{s}}(\mathcal{B}_j^-) \right\}, \quad (3.65)$$

where $\mathcal{B}_j^- = \mathcal{B}_- \cup \Psi_j^-$ with Ψ_j^- defined in (3.5). Then there are three possibilities of how \mathcal{A}_h is going to be updated:

- (i). If $\mathcal{S}_{h,0} = \mathcal{S}_-$, then the forward selection of \mathcal{A}_h will stop and \mathcal{A}_h is removed from \mathcal{V}_+ . For \mathcal{A}_h , its IDF set is $\mathcal{S}_h = \mathcal{S}_{h,0}$.
- (ii). If $\mathcal{S}_{h,0} \subsetneq \mathcal{S}_-$ and $\mathcal{S}_{h,0} \neq \emptyset$, then for any given $\mathbf{s} \in \mathcal{S}_{h,0}$ none of the BICs of the candidate models is lower than the BIC of \mathcal{B}_- , the forward selection of \mathcal{A}_h will stop and \mathcal{A}_h is removed from \mathcal{V}_+ . \mathcal{S}_h is updated by setting $\mathcal{S}_h = \mathcal{S}_{h,0}$. For each $\mathbf{s} \in \mathcal{S}_- \setminus \mathcal{S}_{h,0}$, the procedure finds the set $\Psi_{j^*(\mathbf{s})}^-$ in Ψ^- such that the model $\mathcal{B}_{j^*(\mathbf{s})}^- = \mathcal{B}_- \cup \Psi_{j^*(\mathbf{s})}^-$ has the smallest BIC. The index $j^*(\mathbf{s})$ is obtained using

$$j^*(\mathbf{s}) = \underset{1 \leq j \leq |\Psi^-|}{\text{argmin}} \text{BIC}_{\mathbf{s}}(\mathcal{B}_j^-). \quad (3.66)$$

Note that for $\mathbf{s}_1 \neq \mathbf{s}_2$ where $\mathbf{s}_1, \mathbf{s}_2 \in \mathcal{S}_- \setminus \mathcal{S}_{h,0}$, $j^*(\mathbf{s}_1)$ and $j^*(\mathbf{s}_2)$ could be the same. Denote the unique elements of $\{j^*(\mathbf{s}) : \mathbf{s} \in \mathcal{S}_- \setminus \mathcal{S}_{h,0}\}$ as j_1^*, \dots, j_q^* , then $\Psi_{j_1^*}^-, \dots, \Psi_{j_q^*}^-$ are the unique sets of basis functions selected at the current iteration. For each $\Psi_{j_i^*}^-$ a new model $\mathcal{A}_{h,i}$ is generated:

$$\mathcal{A}_{h,i} = \begin{cases} \mathcal{B}_- \cup \{\psi_{j_i^*}^-\}, & \text{if } j_i^* \leq L_-, \\ \mathcal{B}_- \cup \{\psi_{j_i^*-L_-}^-, \psi_{\lfloor j_i^*-L_- \rfloor}^-\}, & \text{if } j_i^* > L_-, \end{cases} \quad (3.67)$$

for $i = 1, \dots, q$. The IDF set $\mathcal{S}_{h,i}$ for $\mathcal{A}_{h,i}$ is defined as

$$\mathcal{S}_{h,i} = \left\{ \mathbf{s} : \mathbf{s} \in \mathcal{S}_- \setminus \mathcal{S}_{h,0}, \text{BIC}_{\mathbf{s}}(\mathcal{B}_- \cup \Psi_{j_i}^-) = \min_{1 \leq j \leq |\Psi_-|} \text{BIC}_{\mathbf{s}}(\mathcal{B}_j^-) \right\}, \quad (3.68)$$

for $i = 1, \dots, q$. To add these new sets to \mathcal{D} and \mathcal{V}_+ , the procedure generates $\mathcal{A}_{d+i} = \mathcal{A}_{h,i}$ together with $\mathcal{S}_{d+i} = \mathcal{S}_{h,i}$, for $i = 1, \dots, q$, since there are d sets already in \mathcal{D} . $\mathcal{A}_{d+1}, \dots, \mathcal{A}_{d+q}$ are then added to both \mathcal{D} and \mathcal{V}_+ .

- (iii). If $\mathcal{S}_{h,0} = \emptyset$, then the same new sets $\mathcal{A}_{h,1}, \dots, \mathcal{A}_{h,q}$ as in (ii) are generated. $\mathcal{A}_h \in \mathcal{D}$ and $\mathcal{A}_h \in \mathcal{V}_+$ are then updated by setting $\mathcal{A}_h = \mathcal{A}_{h,1}$, and its IDF set is updated by setting $\mathcal{S}_h = \mathcal{S}_{h,1}$. To add the rest of the new models to \mathcal{D} and \mathcal{V}_+ if $q > 1$, the procedure generates $\mathcal{A}_{d+i-1} = \mathcal{A}_{h,i}$ together with $\mathcal{S}_{d+i-1} = \mathcal{S}_{h,i}$, for $i = 2, \dots, q$. $\mathcal{A}_{d+2}, \dots, \mathcal{A}_{d+q}$ are then added to both \mathcal{D} and \mathcal{V}_+ if $q > 1$.

At the next iteration, the procedure finds the next model available in \mathcal{V}_- and assign it as \mathcal{B}_- . Once the update of the last model in \mathcal{V}_- is finished, the procedure updates \mathcal{V}_- by setting $\mathcal{V}_- = \mathcal{V}_+$. Then the forward selection process goes through the same routine again for each model in \mathcal{V}_- . The forward selection for each model in \mathcal{D} and \mathcal{V}_+ stops when the BIC stops decreasing, or when the number of basis functions exceeds M . The whole forward selection process stops when both \mathcal{V}_- and \mathcal{V}_+ are empty sets.

There are two ways the look-ahead procedure prevents overfitting. First, at most M basis functions are allowed to be selected for each model. Second, to prevent having too many models after the forward selection completes, we can also control the maximum number of models to be returned, denoted as M_d . Thus, we have $|\mathcal{D}| \leq M_d$, and this helps to reduce the computation burden, especially when many libraries are used. For a certain iteration, if q new models are generated, but $|\mathcal{D}| + q > M_d$, then only the first $M_d - |\mathcal{D}|$ of them are added to \mathcal{D} and \mathcal{V}_+ , and the rest are discarded. The maximum

number of models that can be generated in the forward selection is $\prod_{l=1}^L |S_l|$. This happens when every IDF tuning vector in \mathcal{S} leads to a unique model. Therefore, when we set $M_d \geq \prod_{l=1}^L |S_l|$, there is no limitation on the number of models. We used this latter setting in all of our simulations.

Now let us look at a specific example that illustrates the forward selection process in the look-ahead procedure. Consider the Blocks-Curves function defined in Table 2.1. We use similar simulation settings as in Section 2.2, with sample size $n = 512$, grid design points from $1/n$ to 1, and the response values generated with $\text{SNR} = 3$.

We use the same libraries as in Table 2.1 for this example, where $\mathcal{L}_0 = \mathcal{U}_1$ consists of the constant and linear terms, $\mathcal{L}_1 = \mathcal{C}_2$ consists of cubic spline representers, and $\mathcal{L}_2 = \mathcal{T}_0$ consists of step functions. The knots used in \mathcal{L}_1 and \mathcal{L}_2 are grid points $\{(j + 3)/n, j = 1, \dots, n - 6\}$, so both \mathcal{L}_1 and \mathcal{L}_2 contain 506 basis functions. In addition, each model allows at most $M = 30$ basis functions to be selected, including both basis functions from \mathcal{L}_0 . Denote the IDF of \mathcal{L}_i as s_i , for $i = 1, 2$. s_1 and s_2 can take values in $S_1 = S_2 = \{1, 2, 3, 4, 5, 6\}$, so in total there are $6^2 = 36$ IDF tuning vectors in the IDF space $\mathcal{S} = S_1 \times S_2$. The maximum number of candidate models M_d is set to be 36, so none of the new models that could be generated are discarded in the forward selection. The IDF of each basis function in \mathcal{L}_0 is fixed at 1. Let ξ_i be the i th basis function in $\mathcal{L}_0 \cup \mathcal{L}_1 \cup \mathcal{L}_2$, so $\mathcal{L}_0 = \{\xi_1, \xi_2\}$, $\mathcal{L}_1 = \{\xi_3, \dots, \xi_{508}\}$, and $\mathcal{L}_2 = \{\xi_{509}, \dots, \xi_{1014}\}$. We now present the details of the steps that the look-ahead procedure took in performing the forward selection based on a single realization of \mathbf{y} . Note that the basis functions selected can change if a different realization of \mathbf{y} were used.

Initialization: Set $\mathcal{D} = \mathcal{V}_- = \mathcal{V}_+ = \{\mathcal{A}_1\}$ and $\mathcal{S}_1 = \mathcal{S}$, where $\mathcal{A}_1 = \{\xi_1, \xi_2\}$.

Example Iterations:

(1). Set $\mathcal{B}_- = \mathcal{A}_1$ and $\mathcal{S}_- = \mathcal{S}_1$. For our \mathbf{y} , the pair $\{\xi_{863}, \xi_{379}\}$ was selected for all

$\mathbf{s} \in \mathcal{S}_-$, so \mathcal{A}_1 is updated by setting $\mathcal{A}_1 = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}\}$. \mathcal{S}_1 stays the same. After the update, $\mathcal{D} = \mathcal{V}_+ = \mathcal{V}_- = \{\mathcal{A}_1\}$.

- (2). Set $\mathcal{B}_- = \mathcal{A}_1$ and $\mathcal{S}_- = \mathcal{S}_1$. Given the IDF set $\mathcal{S}_{1,0} = \{(s_1, s_2) \in \mathcal{S} : s_2 > 3\}$, none of the BICs of the candidate models is lower than $\text{BIC}_{\mathbf{s}}(\mathcal{B}_-)$, so \mathcal{A}_1 completes forward selection for $\mathbf{s} \in \mathcal{S}_{1,0}$, and is removed from \mathcal{V}_+ . \mathcal{S}_1 is updated by setting $\mathcal{S}_1 = \mathcal{S}_{1,0}$.

Given the IDF set $\mathcal{S}_{1,1} = \{(s_1, s_2) \in \mathcal{S} : s_2 \leq 3\}$, the pair $\{\xi_{709}, \xi_{607}\}$ is selected, so a new model \mathcal{A}_2 is generated with $\mathcal{A}_2 = \mathcal{B}_- \cup \{\xi_{709}, \xi_{607}\} = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}, \xi_{709}, \xi_{607}\}$. The IDF set for \mathcal{A}_2 is $\mathcal{S}_2 = \mathcal{S}_{1,1}$. \mathcal{A}_2 is then added to \mathcal{D} and \mathcal{V}_+ . After the update, $\mathcal{D} = \{\mathcal{A}_1, \mathcal{A}_2\}$ and $\mathcal{V}_+ = \{\mathcal{A}_2\}$. \mathcal{V}_- is updated by setting $\mathcal{V}_- = \{\mathcal{A}_2\}$.

- (3). Set $\mathcal{B}_- = \mathcal{A}_2$ and $\mathcal{S}_- = \mathcal{S}_2$. Given the IDF set $\mathcal{S}_{2,1} = \{(s_1, s_2) \in \mathcal{S} : s_1 \leq 3, s_2 \leq 2\}$, the pair $\{\xi_{85}, \xi_{970}\}$ is selected. Given the IDF set $\mathcal{S}_{2,2} = \{(s_1, s_2) \in \mathcal{S} : s_1 \geq 4 \text{ and } s_2 \leq 2, \text{ or } s_1 \geq 5 \text{ and } s_2 = 3\}$, the single basis $\{\xi_{835}\}$ is selected. Given the IDF set $\mathcal{S}_{2,3} = \{(s_1, s_2) \in \mathcal{S} : s_1 \leq 4, s_2 = 3\}$, the pair $\{\xi_{607}, \xi_{85}\}$ is selected. Thus, \mathcal{A}_2 is updated by setting $\mathcal{A}_2 = \mathcal{B}_- \cup \{\xi_{85}, \xi_{970}\} = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}, \xi_{709}, \xi_{607}, \xi_{85}, \xi_{970}\}$, and \mathcal{S}_2 is updated by setting $\mathcal{S}_2 = \mathcal{S}_{2,1}$. Two new sets \mathcal{A}_3 and \mathcal{A}_4 are generated, where $\mathcal{A}_3 = \mathcal{B}_- \cup \{\xi_{835}\} = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}, \xi_{709}, \xi_{607}, \xi_{835}\}$ and $\mathcal{A}_4 = \mathcal{B}_- \cup \{\xi_{607}, \xi_{85}\} = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}, \xi_{709}, \xi_{607}, \xi_{85}\}$. The IDF sets for \mathcal{A}_3 and \mathcal{A}_4 are $\mathcal{S}_3 = \mathcal{S}_{2,2}$ and $\mathcal{S}_4 = \mathcal{S}_{2,3}$ respectively. \mathcal{A}_3 and \mathcal{A}_4 are then added to \mathcal{D} and \mathcal{V}_+ . After the update, $\mathcal{D} = \{\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4\}$ and $\mathcal{V}_+ = \{\mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4\}$. \mathcal{V}_- is updated by setting $\mathcal{V}_- = \{\mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4\}$.

- (4). Set $\mathcal{B}_- = \mathcal{A}_2$ and $\mathcal{S}_- = \mathcal{S}_2$. For all $\mathbf{s} \in \mathcal{S}_-$, none of the BICs of the candidate models is lower than $\text{BIC}_{\mathbf{s}}(\mathcal{B}_-)$, so \mathcal{A}_2 completes forward selection, and is re-

- moved from \mathcal{V}_+ . After the update, $\mathcal{D} = \{\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4\}$, $\mathcal{V}_+ = \{\mathcal{A}_3, \mathcal{A}_4\}$, and $\mathcal{V}_- = \{\mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4\}$.
- (5). Set $\mathcal{B}_- = \mathcal{A}_3$ and $\mathcal{S}_- = \mathcal{S}_3$. Given the IDF set $\mathcal{S}_{3,0} = \{(s_1, s_2) \in \mathcal{S} : s_1 \geq 5, s_2 = 3\}$, none of the BICs of the candidate models is lower than $\text{BIC}_s(\mathcal{B}_-)$, so \mathcal{A}_3 completes forward selection, and is removed from \mathcal{V}_+ . \mathcal{S}_3 is updated by setting $\mathcal{S}_3 = \mathcal{S}_{3,0}$.
- Given the IDF set $\mathcal{S}_{3,1} = \{(s_1, s_2) \in \mathcal{S} : s_1 \geq 4, s_2 = 1\}$, the pair $\{\xi_{875}, \xi_{932}\}$ is selected. Given the IDF set $\mathcal{S}_{3,2} = \{(s_1, s_2) \in \mathcal{S} : s_1 \geq 4, s_2 = 2\}$, the single basis $\{\xi_{875}\}$ is selected. Thus, two new models \mathcal{A}_5 and \mathcal{A}_6 are generated with $\mathcal{A}_5 = \mathcal{B}_- \cup \{\xi_{875}, \xi_{932}\} = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}, \xi_{709}, \xi_{607}, \xi_{835}, \xi_{875}, \xi_{932}\}$ and $\mathcal{A}_6 = \mathcal{B}_- \cup \{\xi_{875}\} = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}, \xi_{709}, \xi_{607}, \xi_{835}, \xi_{875}\}$. Their IDF sets are $\mathcal{S}_5 = \mathcal{S}_{3,1}$, and $\mathcal{S}_6 = \mathcal{S}_{3,2}$. \mathcal{A}_5 and \mathcal{A}_6 are then added to \mathcal{D} and \mathcal{V}_+ . After the update, $\mathcal{D} = \{\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4, \mathcal{A}_5, \mathcal{A}_6\}$, $\mathcal{V}_+ = \{\mathcal{A}_4, \mathcal{A}_5, \mathcal{A}_6\}$, and $\mathcal{V}_- = \{\mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4\}$.
- (6). Set $\mathcal{B}_- = \mathcal{A}_4$ and $\mathcal{S}_- = \mathcal{S}_4$. For all $\mathbf{s} \in \mathcal{S}_-$, none of the BICs of the candidate models is lower than $\text{BIC}_s(\mathcal{B}_-)$, so \mathcal{A}_4 completes forward selection, and is removed from \mathcal{V}_+ . After the update, $\mathcal{D} = \{\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4, \mathcal{A}_5, \mathcal{A}_6\}$ and $\mathcal{V}_+ = \{\mathcal{A}_5, \mathcal{A}_6\}$. \mathcal{V}_- is updated by setting $\mathcal{V}_- = \{\mathcal{A}_5, \mathcal{A}_6\}$.
- (7). Set $\mathcal{B}_- = \mathcal{A}_5$ and $\mathcal{S}_- = \mathcal{S}_5$. For all $\mathbf{s} \in \mathcal{S}_-$, the pair $\{\xi_{984}, \xi_{970}\}$ is selected. Thus, \mathcal{A}_5 is updated by setting $\mathcal{A}_5 = \mathcal{B}_- \cup \{\xi_{984}, \xi_{970}\} = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}, \xi_{709}, \xi_{607}, \xi_{835}, \xi_{875}, \xi_{932}, \xi_{984}, \xi_{970}\}$. \mathcal{S}_5 stays the same. After the update, $\mathcal{D} = \{\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4, \mathcal{A}_5, \mathcal{A}_6\}$, $\mathcal{V}_+ = \{\mathcal{A}_5, \mathcal{A}_6\}$, and $\mathcal{V}_- = \{\mathcal{A}_5, \mathcal{A}_6\}$.
- (8). Set $\mathcal{B}_- = \mathcal{A}_6$ and $\mathcal{S}_- = \mathcal{S}_6$. For all $\mathbf{s} \in \mathcal{S}_-$, none of the BICs of the candidate models is lower than $\text{BIC}_s(\mathcal{B}_-)$, so \mathcal{A}_6 completes forward selection, and is removed from \mathcal{V}_+ . After the update, $\mathcal{D} = \{\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4, \mathcal{A}_5, \mathcal{A}_6\}$ and $\mathcal{V}_+ = \{\mathcal{A}_5\}$. \mathcal{V}_- is updated by setting $\mathcal{V}_- = \{\mathcal{A}_5\}$.

- (9). Set $\mathcal{B}_- = \mathcal{A}_5$ and $\mathcal{S}_- = \mathcal{S}_5$. For all $\mathbf{s} \in \mathcal{S}_-$, the single basis $\{\xi_{665}\}$ is selected. Thus, \mathcal{A}_5 is updated by setting $\mathcal{A}_5 = \mathcal{B}_- \cup \{\xi_{665}\} = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}, \xi_{709}, \xi_{607}, \xi_{835}, \xi_{875}, \xi_{932}, \xi_{984}, \xi_{970}, \xi_{665}\}$. \mathcal{S}_5 stays the same. After the update, $\mathcal{D} = \{\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4, \mathcal{A}_5, \mathcal{A}_6\}$, $\mathcal{V}_+ = \{\mathcal{A}_5\}$, and $\mathcal{V}_- = \{\mathcal{A}_5\}$.
- (10). Set $\mathcal{B}_- = \mathcal{A}_5$ and $\mathcal{S}_- = \mathcal{S}_5$. For all $\mathbf{s} \in \mathcal{S}_-$, none of the BICs of the candidate models is lower than $\text{BIC}_{\mathbf{s}}(\mathcal{B}_-)$, so \mathcal{A}_5 completes forward selection, and is removed from \mathcal{V}_+ . After the update, $\mathcal{D} = \{\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4, \mathcal{A}_5, \mathcal{A}_6\}$ and $\mathcal{V}_+ = \emptyset$. \mathcal{V}_- is updated by setting $\mathcal{V}_- = \emptyset$. Therefore, the whole forward selection process has ended.

We can see from above that whenever new sets are added to \mathcal{D} and \mathcal{V} , the IDF sets for these new models together with those for the updated current model form a partition of the IDF set for \mathcal{B}_- . For example, at iteration (2), $\mathcal{S}_{1,0}$ and $\mathcal{S}_{1,1}$ form a partition of \mathcal{S}_- . The IDF sets for all the sets in \mathcal{D} at iterations (1), (2), (3), and (5) are shown in Figure 3.2. We can clearly see that as the number of models in \mathcal{D} increases, \mathcal{S} is divided into more and more parts, where each part corresponds to a unique model in \mathcal{D} . Also, the boundaries between different parts in \mathcal{S} need not be parallel or perpendicular, as can be seen from iteration (3).

We can also use the tree diagram to summarize the forward selection process described above, as shown in Figure 3.1. Here the edges are labeled with the iteration numbers, and each node represents a unique model with its IDF set. The set \mathcal{V}_- at each round corresponds to the collection of internal nodes with the same depth in the tree. For example, after iteration (3), \mathcal{V}_- becomes $\{\mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4\}$, which can be represented by the collection of all three nodes with depth equal to 3 in the tree diagram. After each model in \mathcal{V}_- is updated at iteration (4) through (6), \mathcal{V}_- is updated by setting $\mathcal{V}_- = \{\mathcal{A}_5, \mathcal{A}_6\}$, since the nodes containing \mathcal{A}_5 and \mathcal{A}_6 are the only two internal nodes at the next level of the tree.

At each iteration, new basis functions that are added to the current model are colored in cyan. For example, from the diagram we can see that at iteration (1), the pair of bases $\{\xi_{863}, \xi_{379}\}$ is selected and added to \mathcal{A}_1 . Each end node is red-filled and corresponds to the end of forward selection for each model. Thus, it can be seen that \mathcal{A}_1 through \mathcal{A}_6 completes the forward selection at iteration (2), (4), (5), (6), (10), and (8) respectively. In addition, when a node has multiple branches, this means the IDF set of the current model is partitioned into multiple parts, where each part leads to a unique model (the child node) after the iteration.

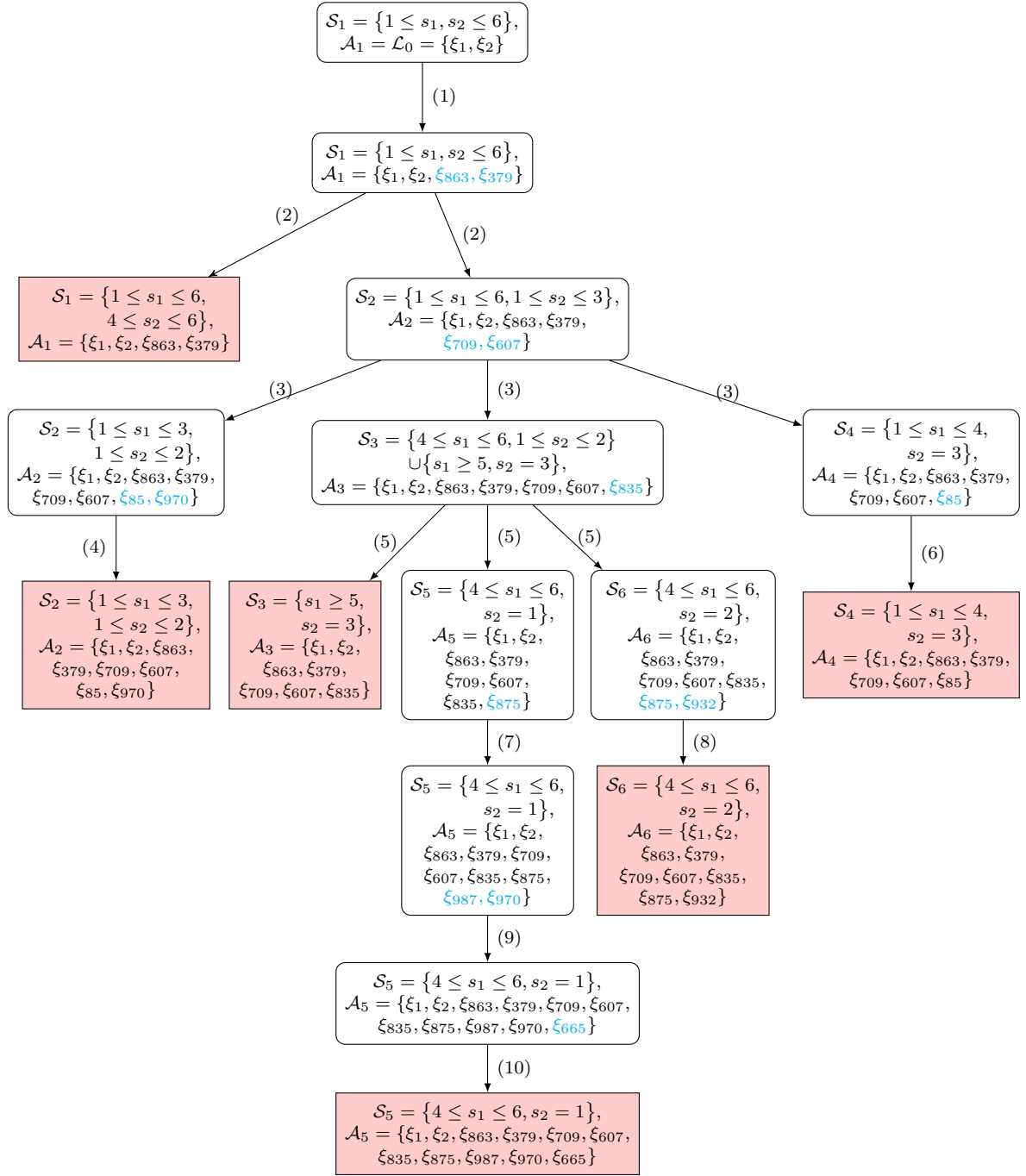
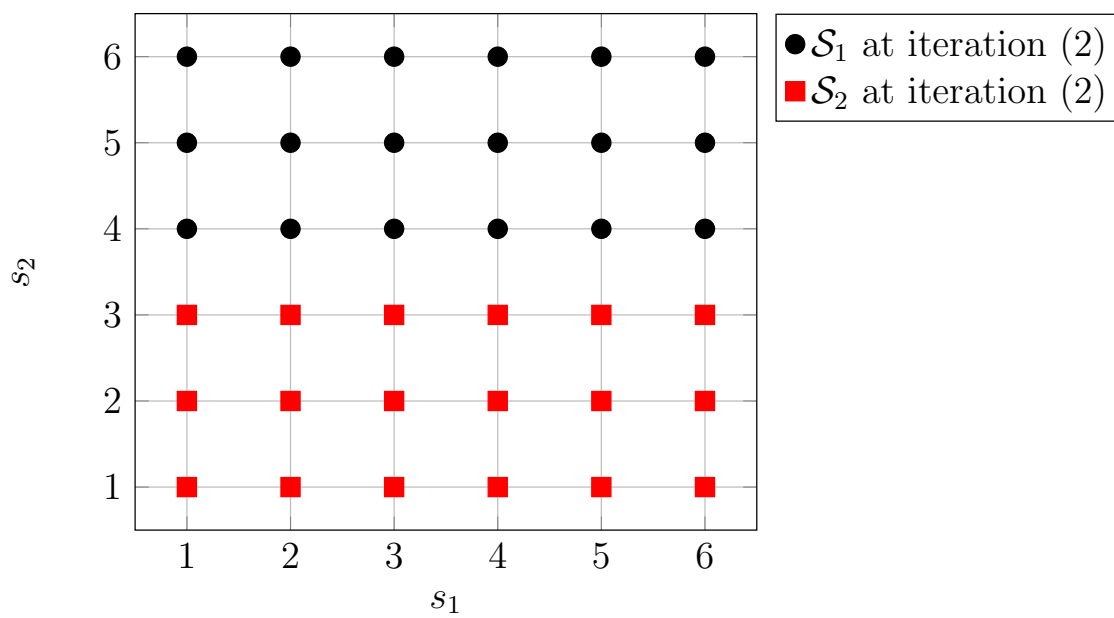
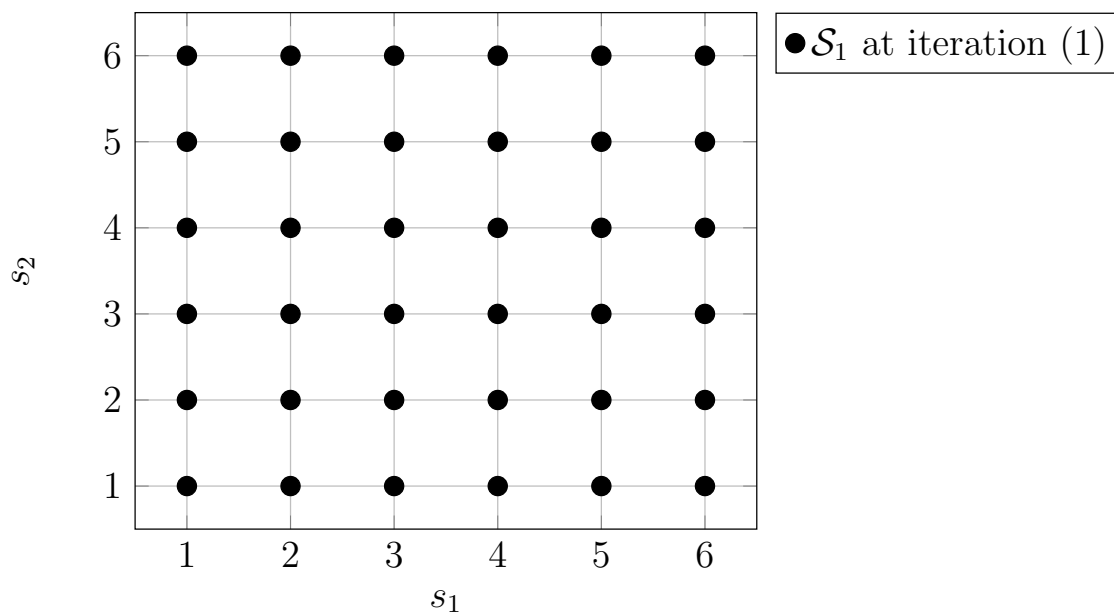


Figure 3.1: Tree diagram of the forward selection process for a data set from the Blocks-Curves example. Note that s_1 and s_2 take integer values from 1 to 6 in this illustration. The numbers in parentheses on left or right side of arrow represent iteration number in the searching procedure. Each end node is red-filled. New basis functions that are selected at each iteration are colored in cyan.



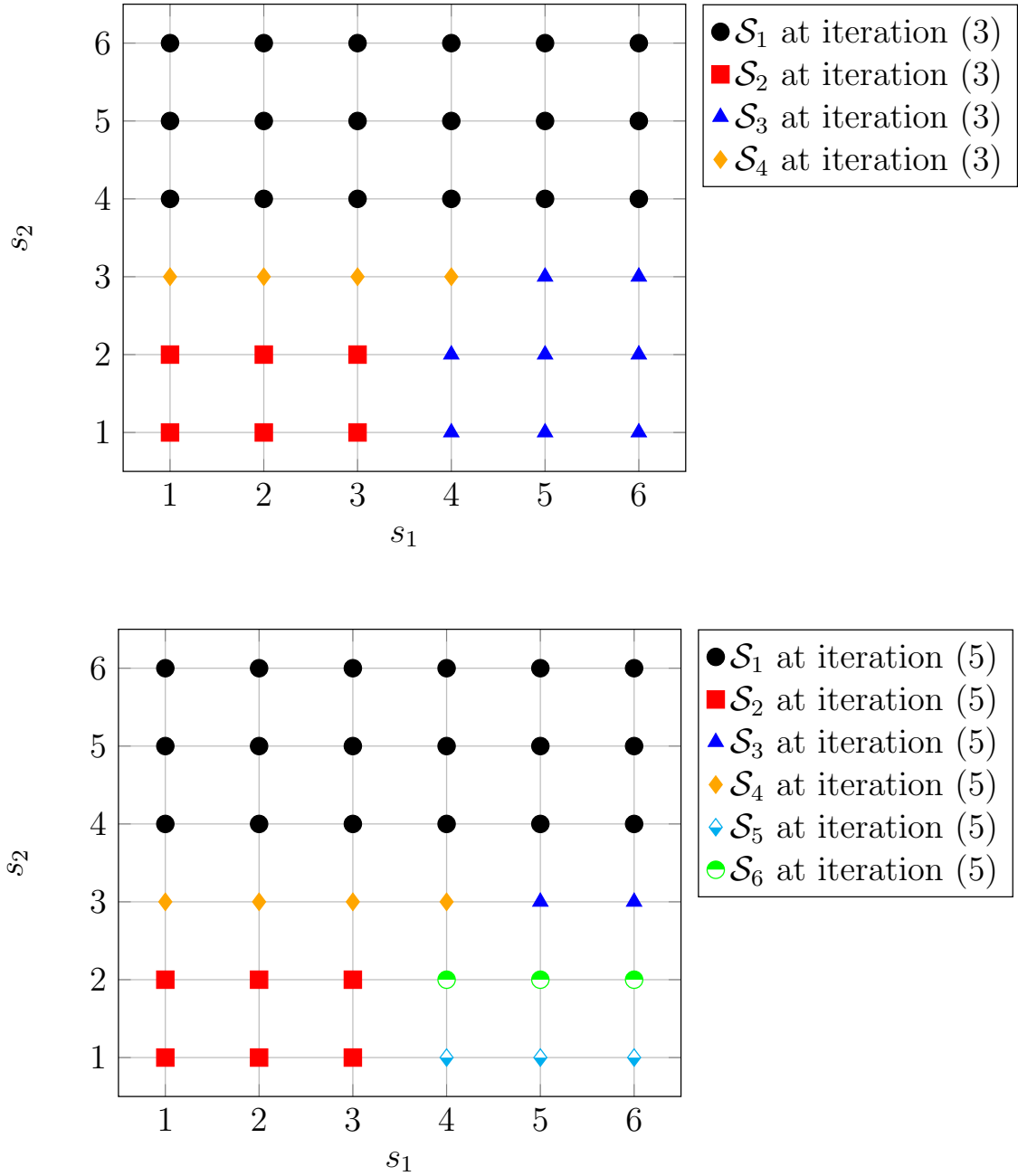


Figure 3.2: The sets of IDF tuning vectors for all the models in \mathcal{D} at different iterations for a data set from the Blocks-Curves example

3.5 Selection of IDF's

After the forward selection process is completed, the look-ahead procedure returns a set of models \mathcal{D} . Each model has its own IDF set which contains at least one IDF tuning vector in the IDF space \mathcal{S} . Clearly if $|\mathcal{D}| < M_d$, then $\{\mathcal{S}_1, \dots, \mathcal{S}_{|\mathcal{D}|}\}$ forms a partition of \mathcal{S} . Otherwise $\{\mathcal{S}_1, \dots, \mathcal{S}_{|\mathcal{D}|}\}$ are disjoint subsets of \mathcal{S} . Now the question is how to select one of the models to be our final model.

Selecting a model in \mathcal{D} is equivalent to selecting one of the IDF tuning vectors in \mathcal{S} , because for a fixed IDF vector, the model is uniquely determined by the forward selection process. The role of the IDF tuning vector is similar to the role of the smoothing parameter λ in the smoothing spline model. By penalizing libraries differently via IDFs, our procedure controls the selection of each type of basis functions. For instance, for the Blocks-Curves example presented in Section 3.4, when the procedure correctly selects the suitable penalties on \mathcal{L}_1 and \mathcal{L}_2 with $s_1 \leq 4$ and $s_2 = 3$, model $\mathcal{A}_4 = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}, \xi_{709}, \xi_{607}, \xi_{85}\}$ will be the final model and exactly three step functions with jump points close to 0.2, 0.4, and 0.7 are selected.

We will use k -fold cross validation to select the best IDF tuning vector in \mathcal{S} . The whole data set is split into k roughly equal-sized parts. For $v = 1, \dots, k$, we use the v -th part as the validation set and the other parts combined as the training set. For a tuning vector $\mathbf{s} \in \mathcal{S}$, define its average cross validation error as:

$$\text{CV}_{\mathbf{s}} = \frac{1}{k} \sum_{v=1}^k \left\{ \frac{1}{n_v} \sum_{i=1}^{n_v} \left(y_i^{[v]} - \hat{f}_{\mathbf{s}}^{[-v]}(\mathbf{x}_i^{[v]}) \right)^2 \right\}, \quad (3.69)$$

where $\{(\mathbf{x}_i^{[v]}, y_i^{[v]}), 1 \leq i \leq n_v\}$ is the v -th validation set, and $\hat{f}_{\mathbf{s}}^{[-v]}$ is the ordinary least squares fit based on all observations except those in the v -th validation set, using the basis functions obtained from the forward selection for the fixed IDF tuning vector \mathbf{s} .

The best IDF tuning vector is chosen to minimize $CV_{\mathbf{s}}$:

$$\mathbf{s}^* = \underset{\mathbf{s} \in \mathcal{S}}{\operatorname{argmin}} CV_{\mathbf{s}}. \quad (3.70)$$

Let \mathcal{S}^* be the IDF set in $\{\mathcal{S}_1, \dots, \mathcal{S}_{|\mathcal{D}|}\}$ that contains \mathbf{s}^* , then the final model \mathcal{A}^* is the model in \mathcal{D} whose IDF set is \mathcal{S}^* .

For the same Blocks-Curves example presented in Section 3.4, here is the levelplot of the average cross validation error for each IDF tuning vector in $\{(s_1, s_2), s_1, s_2 \in \{1, 2, 3, 4, 5, 6\}\}$. We can see that given $s_1 \in \{1, 2, 3\}$ and $s_2 = 3$, the average cross validation error is the smallest. Thus, model $\mathcal{A}_4 = \{\xi_1, \xi_2, \xi_{863}, \xi_{379}, \xi_{709}, \xi_{607}, \xi_{85}\}$ is chosen to be the final model for this specific example.

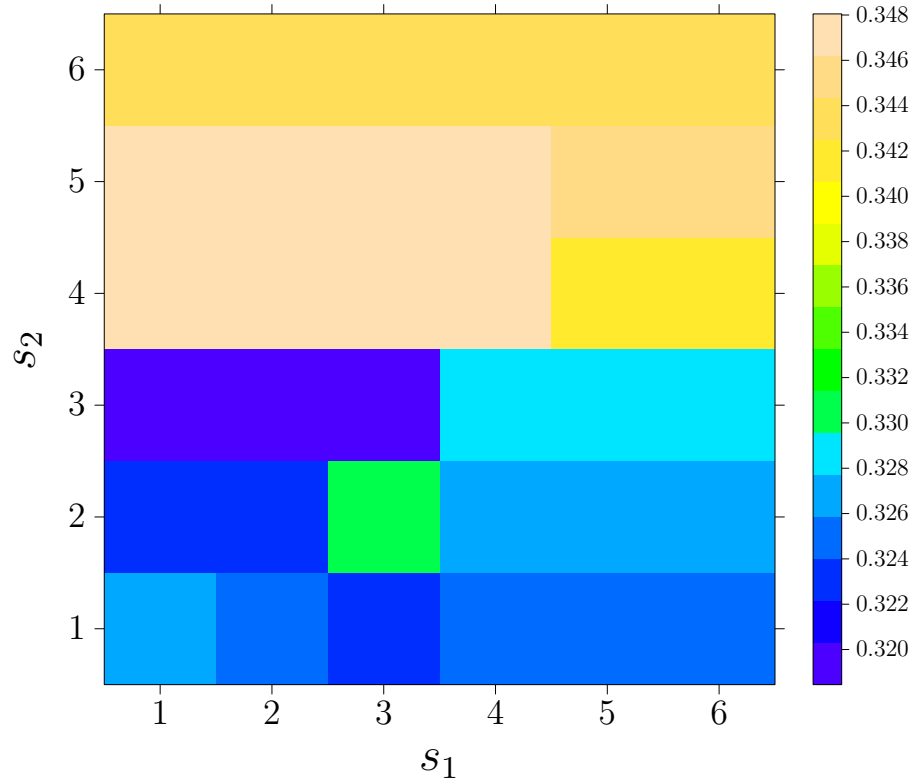


Figure 3.3: The average cross validation error for each IDF tuning vector, using a data set from the Blocks-Curves example.

For the choice of IDF tuning values for each library, one can first try values within a wide range, for example the $\{1, 2, 3, 4, 5, 6\}$ used above. Based on the IDF values selected by LAP, we can then fine-tune the IDFs and apply the LAP again. In the example above, LAP selects $s_1^* \in \{1, 2, 3\}$ and $s_2^* = 3$. For s_1 , we can then specify a finer grid points between 1 and 3, and apply the LAP again to see if the basis functions selected will be different. For this example, the results stay the same.

3.6 The Whole Look-Ahead Procedure

The whole look-ahead procedure is summarized below.

The Look-Ahead Procedure

1. *Initialization:* Set $\mathcal{A}_1 = \mathcal{L}_0$, let M be an upper bound on the number of basis functions to be selected (including those in \mathcal{L}_0), and M_d be the maximum number of models to be generated in the forward selection. Let \mathcal{S} be the IDF space. Set $\mathcal{D} = \mathcal{V}_+ = \mathcal{V}_- = \{\mathcal{A}_1\}$. Set $\mathcal{S}_1 = \mathcal{S}$, $\mathcal{B}_- = \mathcal{A}_1$ and $\mathcal{S}_- = \mathcal{S}_1$.

2. *Forward selection:*

while $\mathcal{V}_+ \neq \emptyset$ and $\mathcal{V}_- \neq \emptyset$ **do**

Find $\mathcal{A}_h \in \mathcal{V}_-$, which the first set in \mathcal{V}_- that has not been updated.

Set $\mathcal{B}_- = \mathcal{A}_h$, and set the IDF set of \mathcal{A}_h as \mathcal{S}_- .

Set $k = |\mathcal{B}_-| + 1$ and $d = |\mathcal{D}|$. Set $\Psi^- = \emptyset$.

for $l = 1, \dots, L$ **do**

if $k = M$ **then**

Select ψ_l^- using (3.1). Update Ψ^- with $\Psi^- \cup \{\psi_l^-\}$.

else

Select ψ_l^- using (3.1). Update Ψ^- with $\Psi^- \cup \{\psi_l^-\}$.

Select $\psi_{\cdot|l}^-$ using (3.2). Update Ψ^- with $\Psi^- \cup \{\psi_l^-, \psi_{\cdot|l}^-\}$.

end if

end for

Let Ψ_j^- be the j th element of Ψ^- . Let $\mathcal{B}_j^- = \mathcal{B}_- \cup \Psi_j^-$ for $j = 1, \dots, |\Psi^-|$.

Set $\mathcal{S}_{h,0} = \emptyset, \Psi_*^- = \emptyset$.

for $s \in \mathcal{S}_-$ **do**

if $\text{BIC}_s(\mathcal{B}_-) \leq \min_{1 \leq j \leq |\Psi^-|} \text{BIC}_s(\mathcal{B}_j^-)$ **then**

 Set $\mathcal{S}_{h,0} = \mathcal{S}_{h,0} \cup \{s\}$.

else

 Solve $j^*(s) = \underset{1 \leq j \leq |\Psi^-|}{\text{argmin}} \text{BIC}_s(\mathcal{B}_j^-)$, with BIC_s defined in (3.6).

 Set $\Psi_*^- = \Psi_*^- \cup \{\Psi_{j^*(s)}^-\}$.

end if

end for

if $\mathcal{S}_{h,0} = \mathcal{S}_-$ **then**

\mathcal{A}_h is removed from \mathcal{V}_+ . The IDF set of \mathcal{A}_h is $\mathcal{S}_h = \mathcal{S}_{h,0}$.

else

 Denote the elements in Ψ_*^- as $\Psi_{j_1^*}^-, \dots, \Psi_{j_q^*}^-$. Generate the new models $\mathcal{A}_{h,1}, \dots, \mathcal{A}_{h,q}$ using (3.67), and their IDF sets $\mathcal{S}_{h,1} \dots \mathcal{S}_{h,q}$, where

$\mathcal{S}_{h,i} = \left\{ s : s \in \mathcal{S}_- \setminus \mathcal{S}_{h,0}, \text{BIC}_s(\mathcal{B}_{j_i^*}^-) = \min_{1 \leq j \leq |\Psi^-|} \text{BIC}_s(\mathcal{B}_j^-) \right\},$

 for $i = 1, \dots, q$.

if $\mathcal{S}_{h,0} \subsetneq \mathcal{S}_-$ and $\mathcal{S}_{h,0} \neq \emptyset$ **then**

 Remove \mathcal{A}_h from \mathcal{V}_+ .

 Generate $\mathcal{A}_{d+i} = \mathcal{A}_{h,i}$ and $\mathcal{S}_{d+i} = \mathcal{S}_{h,i}$ for $i = 1, \dots, q$.

 Add the first $\min(q, M_d - d)$ of $\mathcal{A}_{d+1}, \dots, \mathcal{A}_{d+q}$ to \mathcal{D} and \mathcal{V}_+ .


```

    For  $\mathcal{A}_j \in \mathcal{V}_+$ , remove it from  $\mathcal{V}_+$  if  $|\mathcal{A}_j| = M$ .
  else
    Update  $\mathcal{A}_h$  by setting  $\mathcal{A}_h = \mathcal{A}_{h,1}$ . Update  $\mathcal{S}_h$  by setting  $\mathcal{S}_h = \mathcal{S}_{h,1}$ .
    if  $q > 1$  then
      Generate  $\mathcal{A}_{d+i-1} = \mathcal{A}_{h,i}$  and  $\mathcal{S}_{d+i-1} = \mathcal{S}_{h,i}$  for  $i = 2, \dots, q$ .
      Add the first  $\min(q-1, M_d-d)$  of  $\mathcal{A}_{d+1}, \dots, \mathcal{A}_{d+q-1}$  to  $\mathcal{D}$  and
       $\mathcal{V}_+$ . For  $\mathcal{A}_j \in \mathcal{V}_+$ , remove it from  $\mathcal{V}_+$  if  $|\mathcal{A}_j| = M$ .
    end if
  end if
end if
if all sets in  $\mathcal{V}_-$  have been updated then
  Update  $\mathcal{V}_-$  by setting  $\mathcal{V}_- = \mathcal{V}_+$ .
end if
end while

```

3. *Elimination:* Compute the average cross validation error $\text{CV}_{\mathbf{s}}$ in (3.69) for each $\mathbf{s} \in \mathcal{S}$. Find the \mathbf{s}^* in (3.70). The final model is \mathcal{A}_{h^*} such that $\mathbf{s}^* \in \mathcal{S}_{h^*}$.

3.7 Bootstrap Confidence Intervals

The nonparametric bootstrap confidence intervals have been well studied. The main difficulty in constructing them is due to the bias, which can not be estimated consistently using the bootstrap method. Consider the nonparametric regression model

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2), \quad i = 1, \dots, n, \quad (3.71)$$

where \mathbf{x}_i 's are design points with arbitrary domain \mathcal{X} . As Cummins et al. (2001) has pointed out, by minimizing the average mean squared error (averaged across design points), the bias and variance of the estimate \hat{f} are deliberately balanced, so $\text{Bias}(\hat{f}(\mathbf{x}_i))^2 / \text{Var}(\hat{f}(\mathbf{x}_i))$ converges to a fixed ratio. When this ratio is large, the confidence interval will be consistently centered at the wrong place, causing low coverage. The bias of \hat{f} tends to be big especially at points where f has sharp curvature or sudden jumps, and the coverage probability will be low at these points.

There are two common approaches to overcome this problem. One is by oversmoothing and inducing an explicit or implicit bias estimator, for example see Härdle and Bowman (1988), Härdle and Marron (1991), Hall (1992a), Eubank and Speckman (1993), Sun and Loader (1994), Härdle et al. (1995), and Xia (1998). The other approach is using undersmoothing to reduce the impact of bias, for example see Bjerve et al. (1985), Hall (1992b), Neumann (1995), Neumann and Polzehl (1998), Picard and Tribouley (2000), Claeskens and Keilegom (2003), McMurphy and Politis (2008). See Hall and Horowitz (2013) for a more comprehensive review. In the smoothing splines literature, Wang and Wahba (1995) compared the average coverage probability of various bootstrap confidence intervals, such as the percentile interval (Efron (1982)) and the pivotal interval (Efron (1981)), with Bayesian confidence intervals (Wahba (1983)) for smoothing splines with

Gaussian data, and found they have similar performance. Cummins et al. (2001) constructed pointwise confidence intervals by selecting local smoothing parameters to yield more uniform pointwise coverage.

For the look-ahead procedure, we construct the bootstrap confidence intervals using the undersmoothing approach. First, the LAP is applied to the data $\{(\mathbf{x}_i, y_i), i = 1 \dots, n\}$ to obtain the estimate \hat{f} . Recall that in LAP, criteria such as the BIC and AIC are used for selecting basis functions, and in order to avoid overfitting, the forward selection process for a candidate model stops once BIC or AIC starts to increase. Therefore, the estimate \hat{f} tends to have big bias, especially at places where f has sharp curvature or sudden jumps. To reduce the bias, we use the undersmoothing approach to obtain a “big” model. Sklar et al. (2013) found that for the BSML procedures the cross validation method generally tends to lead to a model that overfits. We will use the same method for the look-ahead procedure.

Let \mathbf{s}^* be the IDF tuning vector that minimizes the average cross validation error and leads to an estimate \hat{f} , as defined in (3.70). The same LAP forward selection routine in Section 3.2 is then applied again with the IDF vector fixed at \mathbf{s}^* . However, this time instead of stopping the forward selection once the BIC or AIC starts to increase, we let it continue until the number of basis functions in the model reaches M . The forward selection starts with $\mathcal{B}^{(1)} = \mathcal{L}_0$ at the first iteration. At the j th iteration, for $j = 2, 3, \dots$, LAP sets $\mathcal{B}_- = \mathcal{B}^{(j-1)}$, then selects the set of basis functions Ψ_{j*}^- by following the same selection routine in (3.1), (3.2), (3.5), (3.6), (3.7), (3.9), (3.12), (3.13), and sets $\mathcal{B}^{(j)} = \mathcal{B}^{(j-1)} \cup \Psi_{j*}^-$. This process continues until it reaches iteration D such that $|\mathcal{B}^{(D)}| = M$. Therefore, we have D candidate models: $\mathcal{B}^{(1)}, \dots, \mathcal{B}^{(D)}$. Then the k -fold cross validation method is used to select one of them as the final model. The average

cross validation error for model $\mathcal{B}^{(j)}$, $j = 1, \dots, D$, is defined as:

$$\text{CV}(\mathcal{B}^{(j)}) = \frac{1}{k} \sum_{v=1}^k \left\{ \frac{1}{n_v} \sum_{i=1}^{n_v} \left(y_i^{[v]} - \hat{f}_{(j)}^{[-v]}(\mathbf{x}_i^{[v]}) \right)^2 \right\}, \quad (3.72)$$

where $\{(\mathbf{x}_i^{[v]}, y_i^{[v]}), 1 \leq i \leq n_v\}$ is the v -th validation set, and $\hat{f}_{(j)}^{[-v]}$ is the ordinary least squares fit based on all observations except those in the v -th validation set, using the basis functions contained in $\mathcal{B}^{(j)}$. The final model $\tilde{\mathcal{B}}$ is chosen to minimize the average k -fold cross validation error:

$$\tilde{\mathcal{B}} = \underset{\mathcal{B}^{(j)} \in \{\mathcal{B}^{(1)}, \dots, \mathcal{B}^{(D)}\}}{\text{argmin}} \text{CV}(\mathcal{B}^{(j)}). \quad (3.73)$$

Once $\tilde{\mathcal{B}}$ has been chosen, denote \tilde{f} as the ordinary least squares fit using the basis functions contained in $\tilde{\mathcal{B}}$. Then B bootstrap samples $\tilde{\mathbf{y}}_b = (\tilde{y}_{1,b}, \dots, \tilde{y}_{n,b})'$ for $b = 1, \dots, B$ are generated according to

$$\tilde{y}_{i,b} = \tilde{f}(\mathbf{x}_i) + \tilde{\varepsilon}_{i,b}, \quad \tilde{\varepsilon}_{i,b} \stackrel{iid}{\sim} N(0, \hat{\sigma}^2), \quad i = 1, \dots, n; \quad b = 1, \dots, B, \quad (3.74)$$

where $\hat{\sigma}^2$ is an estimate of σ^2 . Typically a difference-based estimator of σ^2 is used, for example see Rice (1984), Gasser et al. (1986), and Tong and Wang (2005). For the look-ahead procedure, the Rice estimator defined in (2.13) is used as the default. Alternatively, the bootstrap samples can be generated by resampling the residuals with replacement. In the look-ahead procedure, only the parametric bootstrap given in (3.74) is implemented.

For each bootstrap sample $\tilde{\mathbf{y}}_b$, $b = 1, \dots, B$, the same forward selection routine mentioned above is applied to obtain the sequence of candidate models $\mathcal{B}_b^{(1)}, \dots, \mathcal{B}_b^{(D_b)}$. To select one of them as the final model, we use the average square error (ASE) instead of the average cross validation error as the model elimination criterion, because the true

function \tilde{f} in the generation of the bootstrap samples is known. The ASE of the model $\mathcal{B}_b^{(j)}$ is defined as:

$$\text{ASE}(\mathcal{B}_b^{(j)}) = \frac{1}{n} \sum_{i=1}^n \left(\hat{f}_b^{(j)}(\mathbf{x}_i) - \tilde{f}(\mathbf{x}_i) \right)^2, \quad (3.75)$$

for $j = 1, \dots, D_b$, where $\hat{f}_b^{(j)}$ is the ordinary least squares fit using the basis functions contained in $\mathcal{B}_b^{(j)}$. For the b th bootstrap sample, the final model \mathcal{B}_b^* is selected to minimize the ASE, i.e.,

$$\mathcal{B}_b^* = \underset{\mathcal{B}_b^{(j)} \in \{\mathcal{B}_b^{(1)}, \dots, \mathcal{B}_b^{(D)}\}}{\text{argmin}} \text{ASE}(\mathcal{B}_b^{(j)}). \quad (3.76)$$

Denote \hat{f}_b^* as the ordinary least squares fit using the basis functions contained in \mathcal{B}_b^* , for $b = 1, \dots, B$. Then the $(1 - \alpha)\%$ confidence interval of $f(\mathbf{x})$ is $(\hat{f}_{\alpha/2}, \hat{f}_{(1-\alpha/2)})$, where $\hat{f}_{\alpha/2}$ and $\hat{f}_{(1-\alpha/2)}$ are the $\alpha/2$ and $(1 - \alpha/2)$ percentiles of $(\hat{f}_1^*(\mathbf{x}), \dots, \hat{f}_B^*(\mathbf{x}))$.

The simulation results for the bootstrap confidence intervals are presented in Section 4.2.

3.8 R Functions for the Look-Ahead Procedure

The look-ahead procedure has been implemented in R. Here we briefly introduce the input arguments and output values of the R functions that can be used for applying LAP and performing predictions.

3.8.1 LAP.R

Description

This function implements the look-ahead procedure (LAP). The program adaptively select one or a pair of basis functions from different libraries according to certain criterion such as the BIC. Each library is associated with an IDF, which is treated as a tuning parameter. K -fold cross validation is used for selecting the best IDF values and the final model.

Usage

```
LAP(y, baseslist, maxbas=30, sub.maxbas=c(), idfs, maxmod=nrow(idfs),
    nfold=10, criterion="BIC")
```

Arguments

- **y**: response vector.
- **baseslist**: list of libraries of basis functions, where each element corresponds on the list corresponds to a library. The first library is the null library \mathcal{L}_0 , the second library is \mathcal{L}_1 and so on. Each library should be a matrix, where each column corresponds to a basis function evaluated at the design points.
- **maxbas**: maximum number of bases allowed to be selected. This is the M in the LAP algorithm.

- `sub.maxbas`: a vector with length equal to the number of nonnull libraries. Each element specifies the maximum number of bases to be preselected for each nonnull library.
- `idfs`: matrix of IDF tuning values with $L + 1$ columns, where each row represents an IDF tuning vector. The IDF for \mathcal{L}_0 (the first column) should be fixed at 1.
- `maxmod`: maximum number of models allowed to be generated during the forward selection process. This is denoted as M_d in the LAP algorithm. The default is the total number of IDF tuning vectors, or equivalently, the number of rows in the matrix `idfs`.
- `nfolds`: the number of folds used in the cross validation.
- `criterion`: basis selection criterion. Available options are “AIC” and “BIC”. Default is “BIC”.

Value

- `bases.chosen`: a matrix with two columns. The second column contains indices of the basis functions selected by LAP. The first column contains the indices of the libraries to which basis functions in the second column belong, where “1” represents the null library, “2” represents \mathcal{L}_1 , and so on.
- `cv.error`: the average cross validation error for each IDF tuning vector, see (3.72).
- `idfs.cv`: the IDF values that minimize average cross validation error, see (3.70).
- `fit`: the fitted values of the final model.
- `coef`: the coefficients in the final model.
- `sigma_hat`: estimated error standard deviation.

- **RSS**: residual sum of squares of the final model.
- **y**: response vector. Same as input.
- **baseslist**: list of libraries of basis functions. Same as input.
- **maxbas**: maximum number of bases allowed to be selected. Same as input.
- **idfs**: list of IDF tuning values for each library, including \mathcal{L}_0 . Same as input.
- **maxmod**: maximum number of models allowed to be generated during the forward selection process. Same as input.
- **criterion**: basis selection criterion. Same as input.

3.8.2 **predict.LAP.R**

Description

This function is for making predictions using the output object from LAP.

Usage

```
predict.LAP (object, bases.include=NULL, new.baseslist=NULL,  
             confint=F, alpha=0.05, sigma=NULL, bootrep=200)
```

Arguments

- **object**: output object from LAP.
- **bases.include**: a vector of length (number of null bases + number of nonnull libraries). The elements should be either 0 or 1, with 0 means the corresponding basis or library will not be included in the prediction, and 1 otherwise.
- **new.baseslist**: the list of libraries used in the prediction. If the original list in the LAP object is used, the prediction will be same as the fitted values in the LAP object.

- **confint**: True or False. If True, a confidence interval with significance level α will be computed. The default is False.
- **alpha**: the significance level for the confidence interval. Default is 0.05.
- **sigma**: the error standard deviation used in constructing the bootstrap sample. The Rice estimator is used as the default.
- **bootrep**: number of bootstrap samples used in constructing the confidence interval. Default is 200.

Value

- **fit**: predicted values.
- **fit.boot**: the fit of the model $\tilde{\mathcal{B}}$ obtained by using cross validation.
- **bases.boot**: the indices of the basis functions contained in the model $\tilde{\mathcal{B}}$.
- **bootfits**: a matrix which contains all the bootstrap fitted values used in constructing the confidence interval as its columns. NULL if **confint**=F.
- **boot.y**: a matrix which contains the bootstrap samples generated for the response variable as its columns. NULL if **confint**=F.
- **bootbases.all**: a list which contains the indices of all the basis functions selected for each bootstrap sample. NULL if **confint**=F.
- **bootbases.pred**: a list which contains the indices of selected basis functions, among all the wanted ones as specified in **bases.include**, for each bootstrap sample. NULL if **confint**=F.
- **lower**: the lower bound of the confidence interval. NULL if **confint**=F.

- **upper**: the upper bound of the confidence interval. NULL if **confint=F**.
- **sigma**: the error standard deviation used in constructing the bootstrap sample.

3.8.3 stdz.R

Description

This function is used to standardize the basis functions (columns) in a design matrix, so that each basis functions evaluated at the design points has norm one.

Usage

```
stdz(X)
```

Arguments

- **X**: matrix of basis functions. Each column of X is a basis function evaluated at the design points.

Value

- matrix with standardized basis functions as its columns.

3.8.4 An Example

We use one simulated data set from the Block-Curves example in Section 3.4 to illustrate the usage of these R functions. The R code is shown below, where the R file `LAP.R` contains `LAP` and `predict.LAP` functions. Here the `stdz` function is used to numerically stabilize the computation. The input matrix in the `stdz` function is standardized by column so that each column has norm one.

```
library(assist)
source("LAP.R")
```

```
set.seed(82)
```

```

n <- 512
x <- seq(from=1/n, to=1, by=1/n)
fx <- 1*(x-0.2 >= 0 & x-0.4 < 0) - x^10*(x-0.7 >= 0) +
      exp((x+0.5)^2)*(x-0.4 >= 0 & x-0.7 < 0)

## Libraries of basis functions:
loc <- x[-c(1,2,3,n-2,n-1,n)]
pt <- rep(1,n)%o%loc
L0 <- stdz(cbind(1, x))
L1 <- stdz((cubic(s = x, t = loc)))
L2 <- stdz((1*(x-pt)>0))
baseslist <- list(L0, L1, L2)

## Look-Ahead Procedure:
y <- fx + rnorm(n, mean=0, sd=sd(fx)/3)

nlib <- length(baseslist)
idfs <- append(list(1),rep(list(c(1:6)),nlib-1))
idfs <- as.matrix(expand.grid(idfs))
fit.lap <- LAP(y, baseslist, maxbas=30, idfs=idfs, criterion="BIC")

## 95% bootstrap confidence interval:
pred.lap <- predict.LAP(fit.lap, bases.include=c(1,1,1,1),
                        new.baseslist=baseslist, confint=T, bootrep=200)

## Plot of the LAP fit with its 95% confidence interval:
plot(x, y, col="gray")
lines(x, fx, type="l", lwd=2)
lines(x, fit.lap$fit, col=2, lty=1, lwd=2)
lines(x, pred.lap$lower, col=2, lty=2, lwd=2)
lines(x, pred.lap$upper, col=2, lty=2, lwd=2)

```

Figure 3.4 shows the LAP fit with its 95% bootstrap confidence interval for this example.

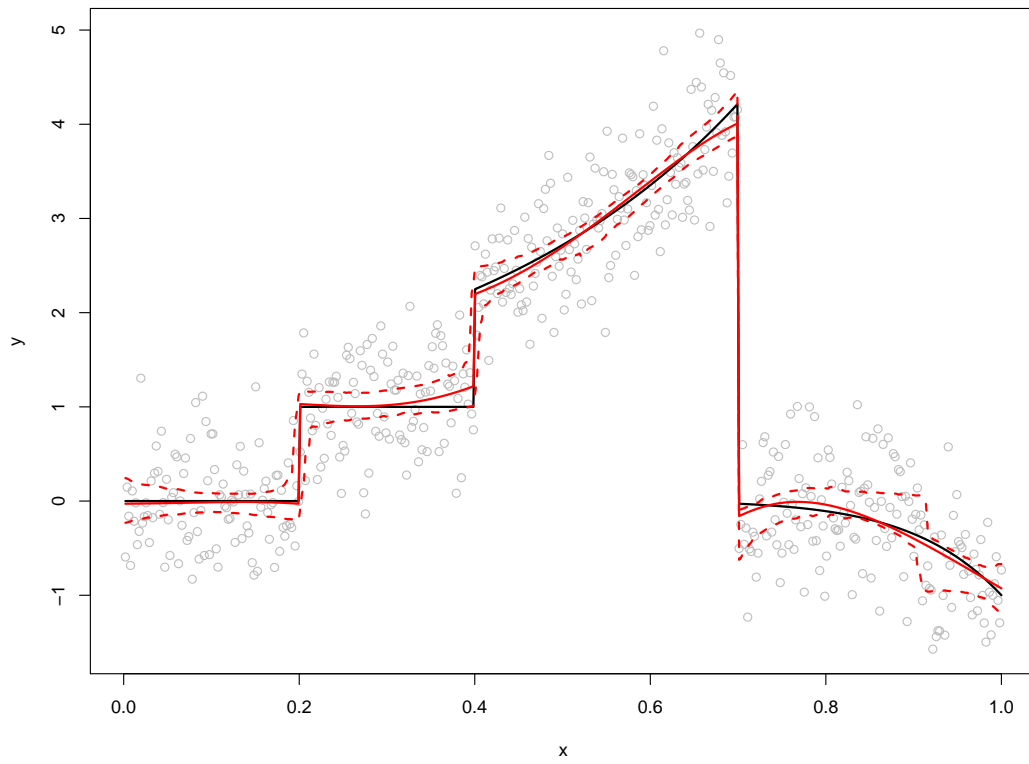


Figure 3.4: LAP fit of the Block-Curves example in Section 3.4. Gray dots are the simulated data points. Black solid line is true function. Red solid line is the LAP fit. Red dashed lines are the upper and lower bounds of the 95% bootstrap confidence interval based on 200 bootstrap samples.

Chapter 4

Simulations

4.1 Estimation of Univariate Functions

We performed simulations to compare the performances of the BSML procedures and our look-ahead procedure with the same six univariate examples presented in Section 2.2. See Table 2.1 for the true functions and the libraries of basis functions used for each example.

The design points used for all simulations are the grid points $\{x_i = i/n : i = 1, \dots, n\}$, where the sample size $n \in \{256, 512, 1024\}$. The response variable is generated using (2.6), where σ is chosen such that $\text{SNR} \in \{1, 2, 3, 4, 5, 6\}$. The knots used in the libraries $\mathcal{P}_2, \mathcal{C}_2, \mathcal{T}_0, \mathcal{T}_2$ are grid points $\{j/n, j = 4, \dots, n - 3\}$. For all procedures, the maximal number of bases M is fixed at 30. For the look-ahead procedure, we used IDF spaces $S_i = \{1, 2, 3, 4, 5, 6\}$ for $i = 1, 2$. The maximal number of models can be returned after the forward selection, denoted by M_d , was set to be $|S_1| \times |S_2| = 36$, so that all new models generated during the forward selection are kept in \mathcal{D} . The 10-fold cross validation is used to select the final model in the look-ahead procedure. For each sample size, 100 random samples are simulated to evaluate the performances of the BSML procedures

and the look-ahead procedure. We use the MSE defined in (2.8) as the measure of performance. For convenience, we will use the acronym “LAP” when referring to the look-ahead procedure from now on. The simulation results are shown in Figure 4.1 – 4.9.

We can see that LAP’s performance is comparable to that of BSML-C and BSML-S in terms of the MSE. For the Sine-Jumps and LW6 examples, on average LAP selected fewer basis function than BSML-C and BSML-S. Also, as we mentioned in Section 2.4, for the Heavisine example the basis functions from \mathcal{L}_1 and \mathcal{L}_2 need to be treated separately. Because of this, for the Heavisine example the performance of BSML-C is not as good as that of BSML-S and LAP across all SNR values and sample sizes. In terms of the computation speed, LAP has clear advantage over both BSML procedures. Table 4.1 – 4.3 lists the average CPU time of fitting each univariate example using BSML-C, BSML-S, and LAP for different samples sizes with $\text{SNR} = 3$. Without having to estimate the GDF at each forward selection step as BSML-C and BSML-S do, LAP manages to be about twice as fast as BSML-C, and more than 10 times faster than BSML-S.

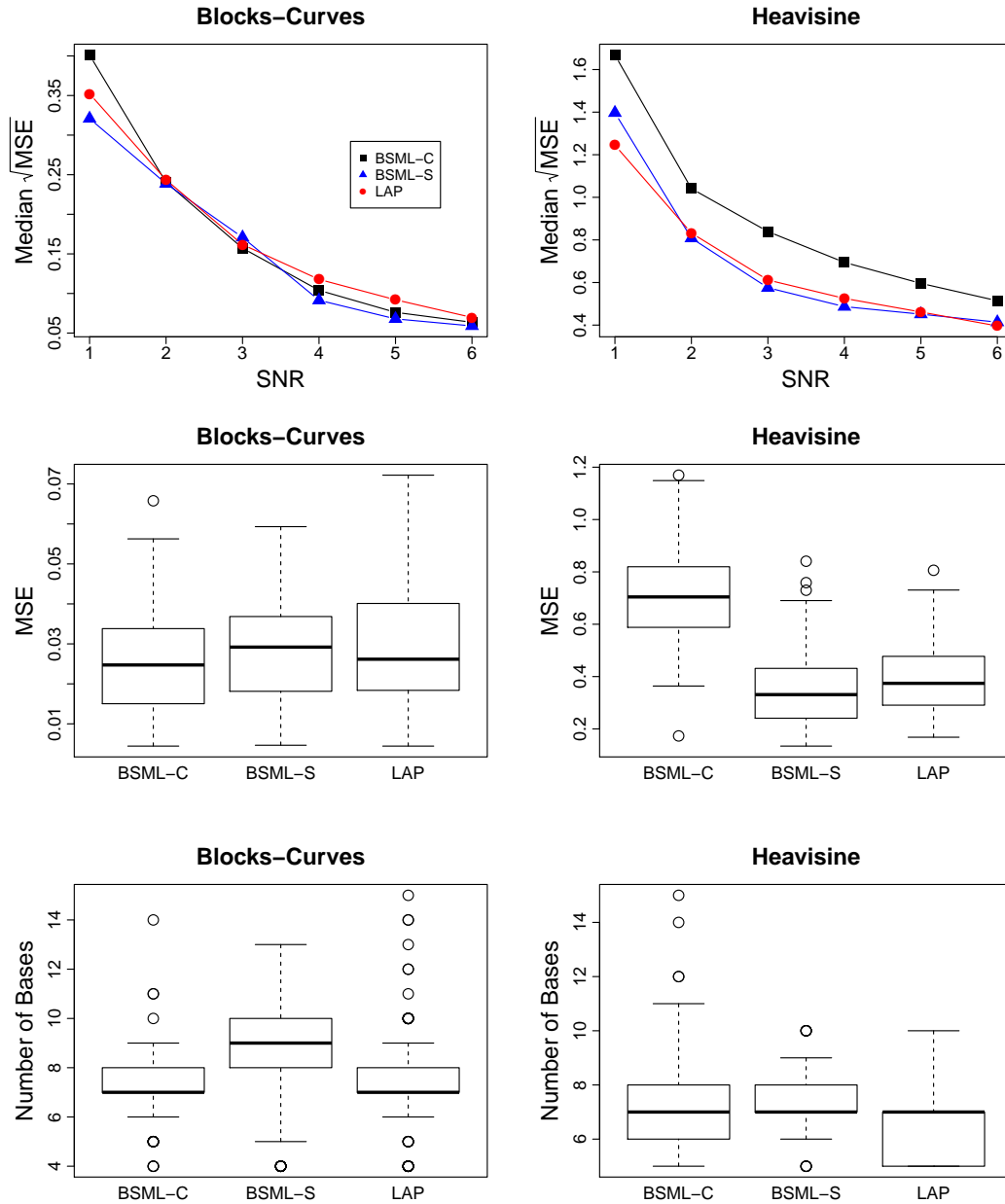


Figure 4.1: Simulation results for the Blocks-Curves and Heavisine examples with $n = 256$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for both examples, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, and \bullet for the look-ahead procedure. The middle panel shows the boxplots of the MSE when SNR = 3. The bottom panel shows the number of basis functions selected when SNR = 3.

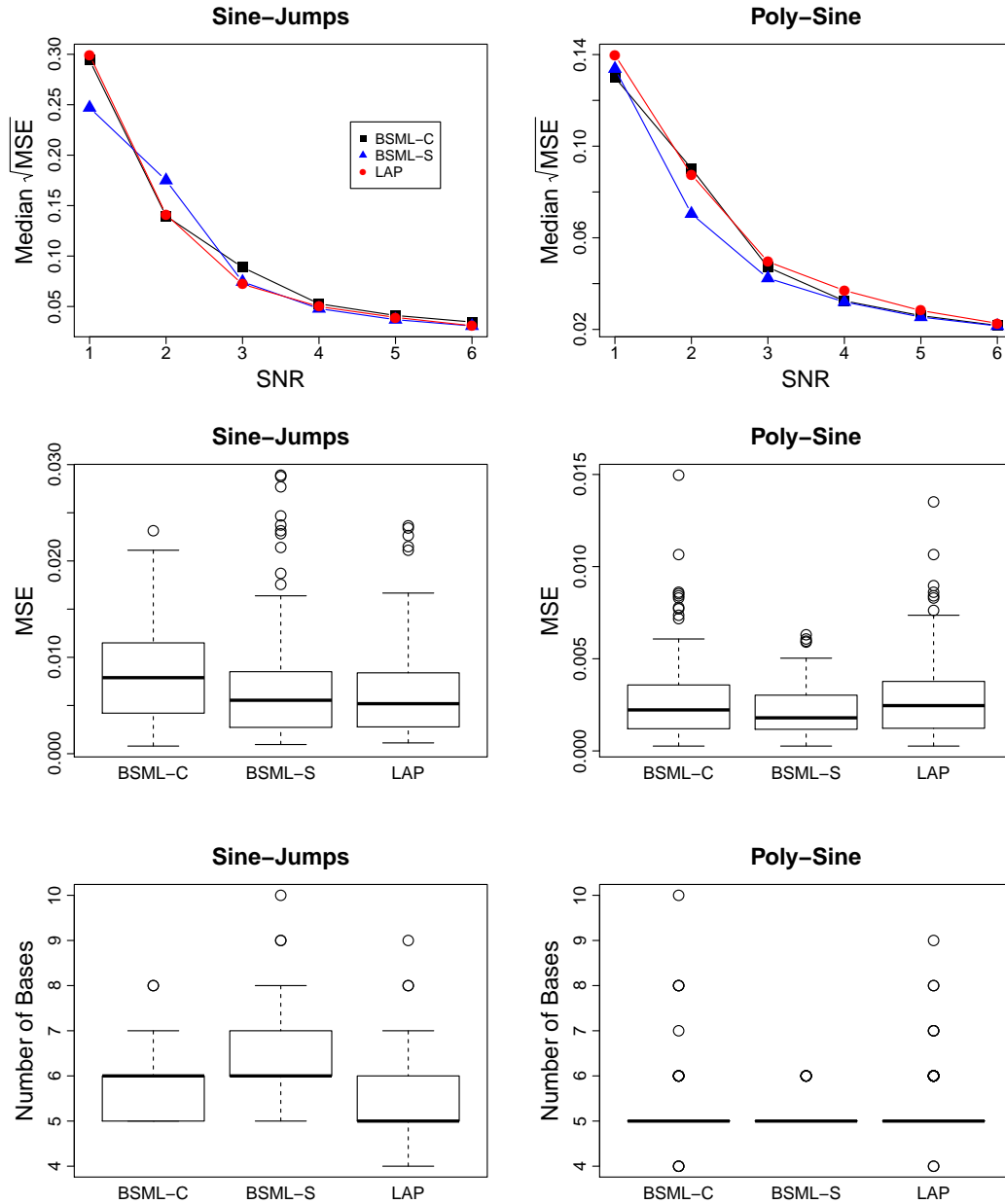


Figure 4.2: Simulation results for the Sine-Jumps and Poly-Sine examples with $n = 256$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for both examples, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, and \bullet for the look-ahead procedure. The middle panel shows the boxplots of the MSE when SNR = 3. The bottom panel shows the number of basis functions selected when SNR = 3.

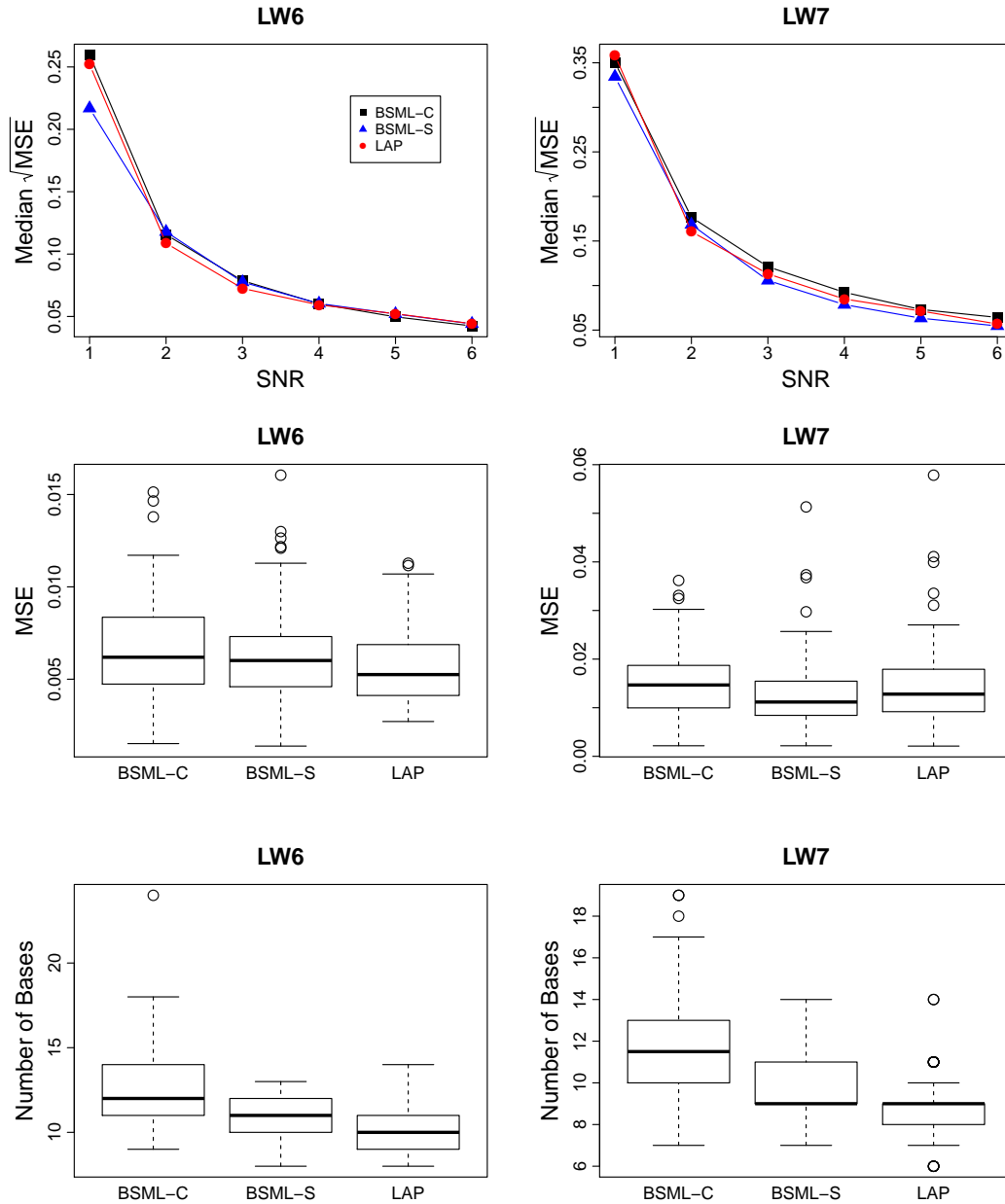


Figure 4.3: Simulation results for the LW6 and LW7 examples with $n = 256$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for both examples, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, and \bullet for the look-ahead procedure. The middle panel shows the boxplots of the MSE when SNR = 3. The bottom panel shows the number of basis functions selected when SNR = 3.

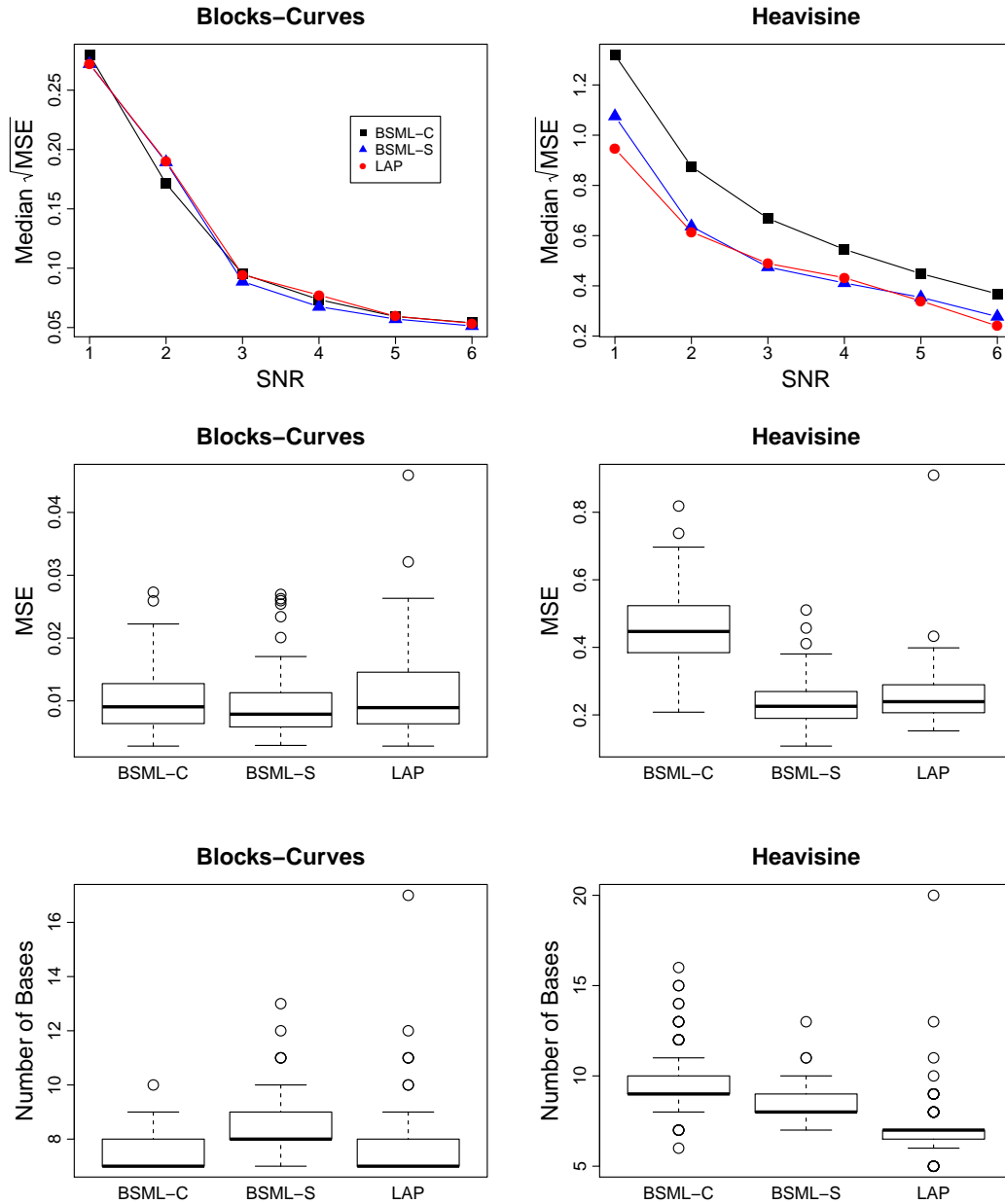


Figure 4.4: Simulation results for the Blocks-Curves and Heavisine examples with $n = 512$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for both examples, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, and \bullet for the look-ahead procedure. The middle panel shows the boxplots of the MSE when SNR = 3. The bottom panel shows the number of basis functions selected when SNR = 3.

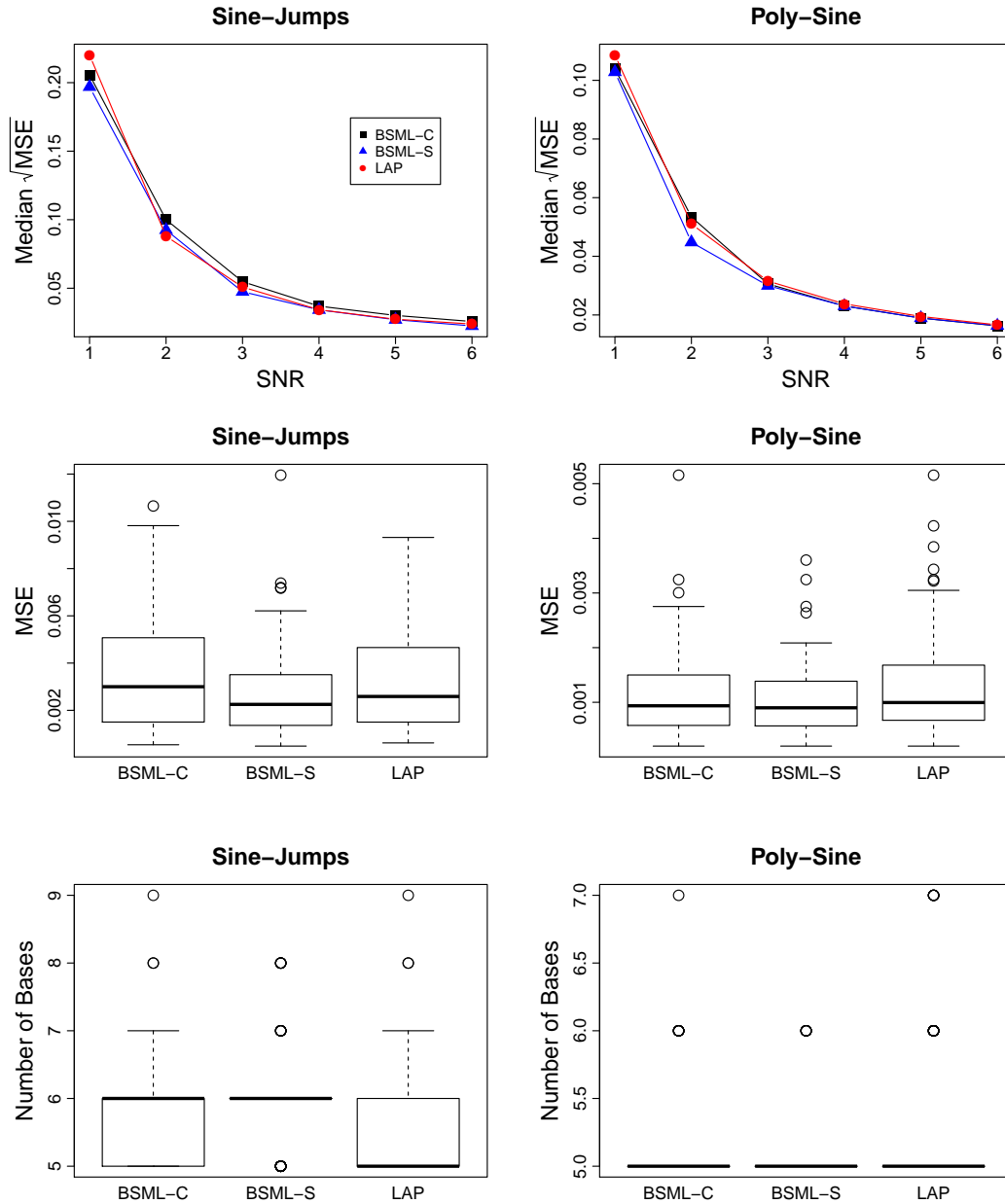


Figure 4.5: Simulation results for the Sine-Jumps and Poly-Sine examples with $n = 512$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for both examples, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, and \bullet for the look-ahead procedure. The middle panel shows the boxplots of the MSE when SNR = 3. The bottom panel shows the number of basis functions selected when SNR = 3.

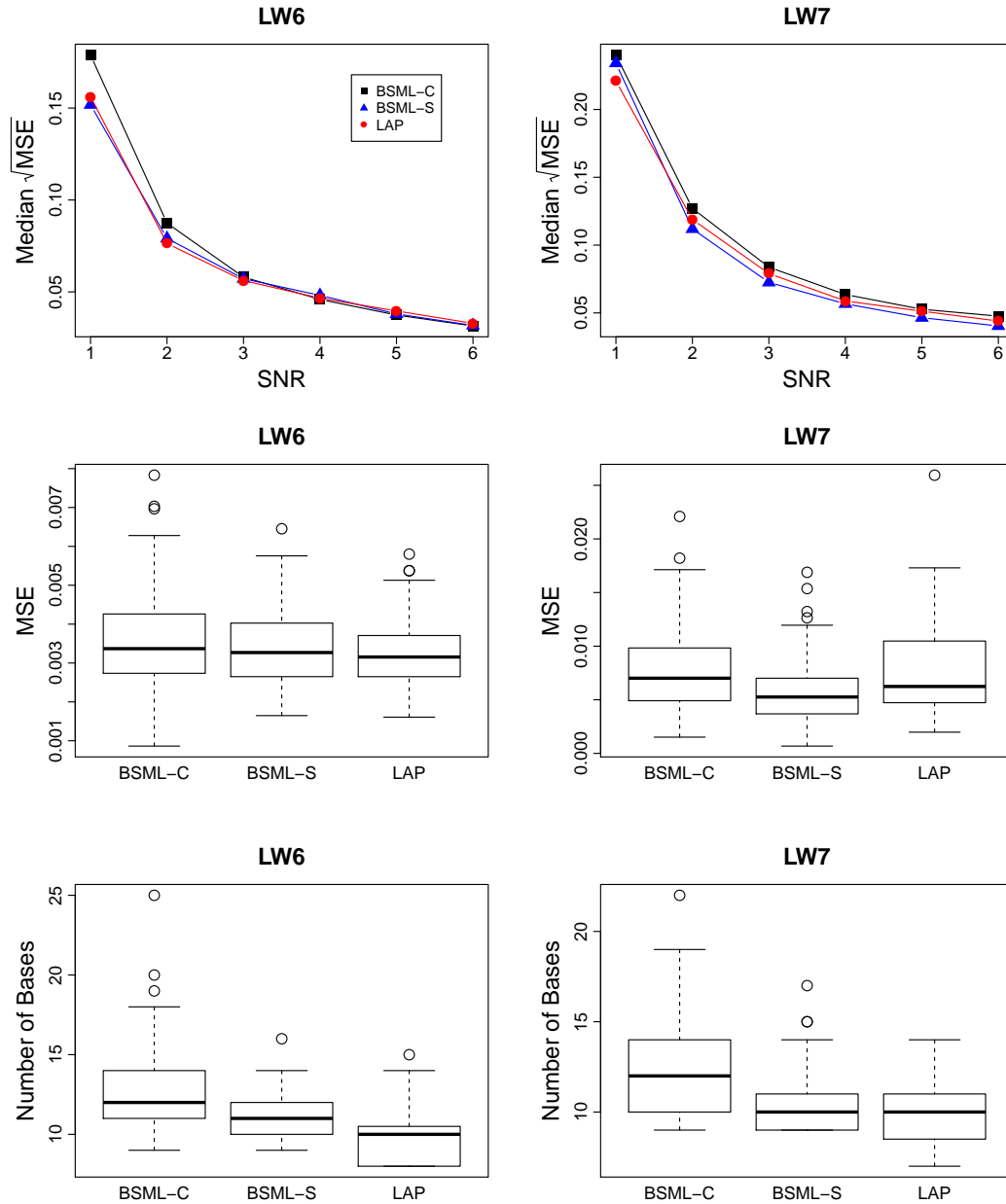


Figure 4.6: Simulation results for the LW6 and LW7 examples with $n = 512$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for both examples, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, and \bullet for the look-ahead procedure. The middle panel shows the boxplots of the MSE when SNR = 3. The bottom panel shows the number of basis functions selected when SNR = 3.

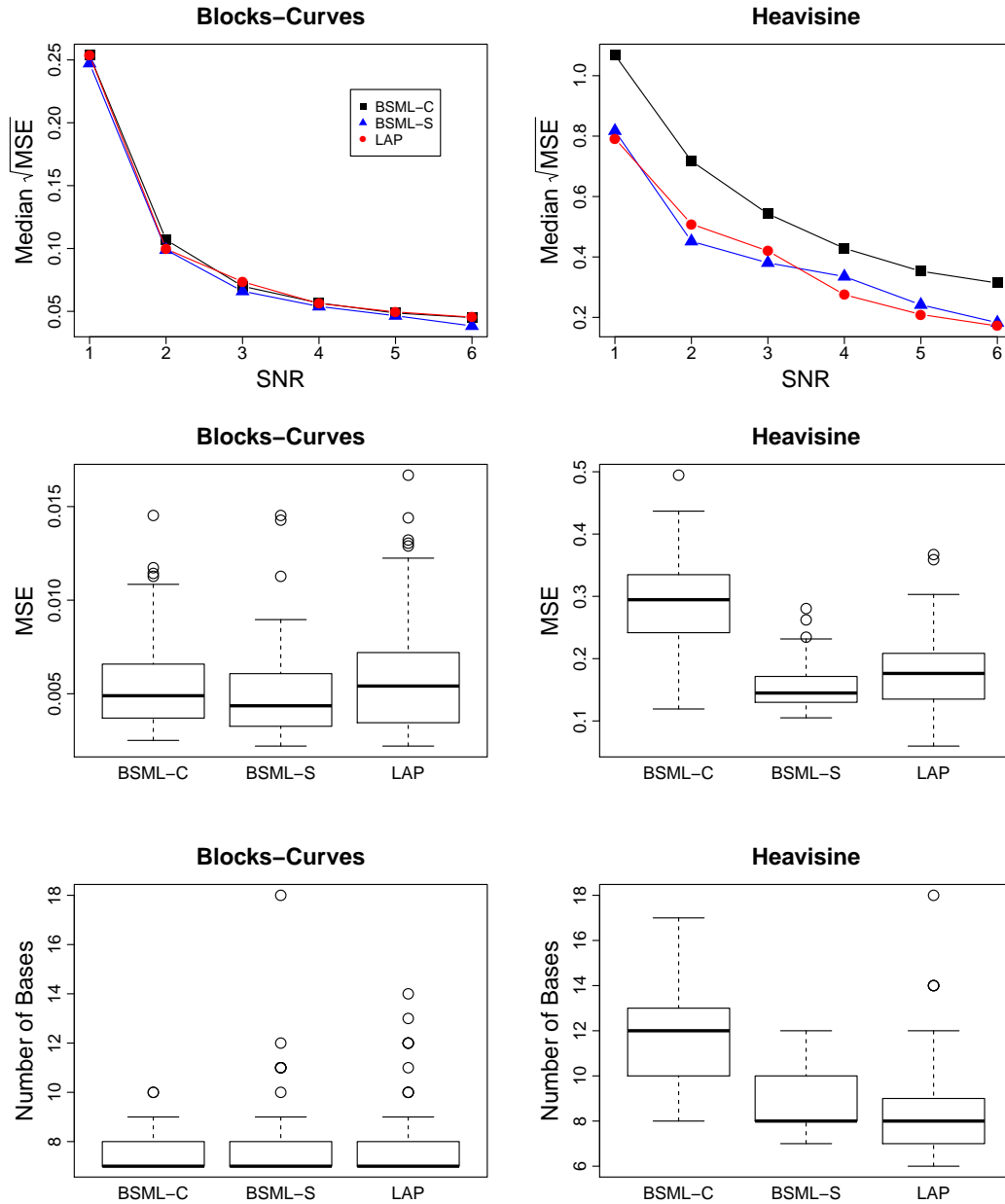


Figure 4.7: Simulation results for the Blocks-Curves and Heavisine examples with $n = 1024$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for both examples, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, and \bullet for the look-ahead procedure. The middle panel shows the boxplots of the MSE when SNR = 3. The bottom panel shows the number of basis functions selected when SNR = 3.

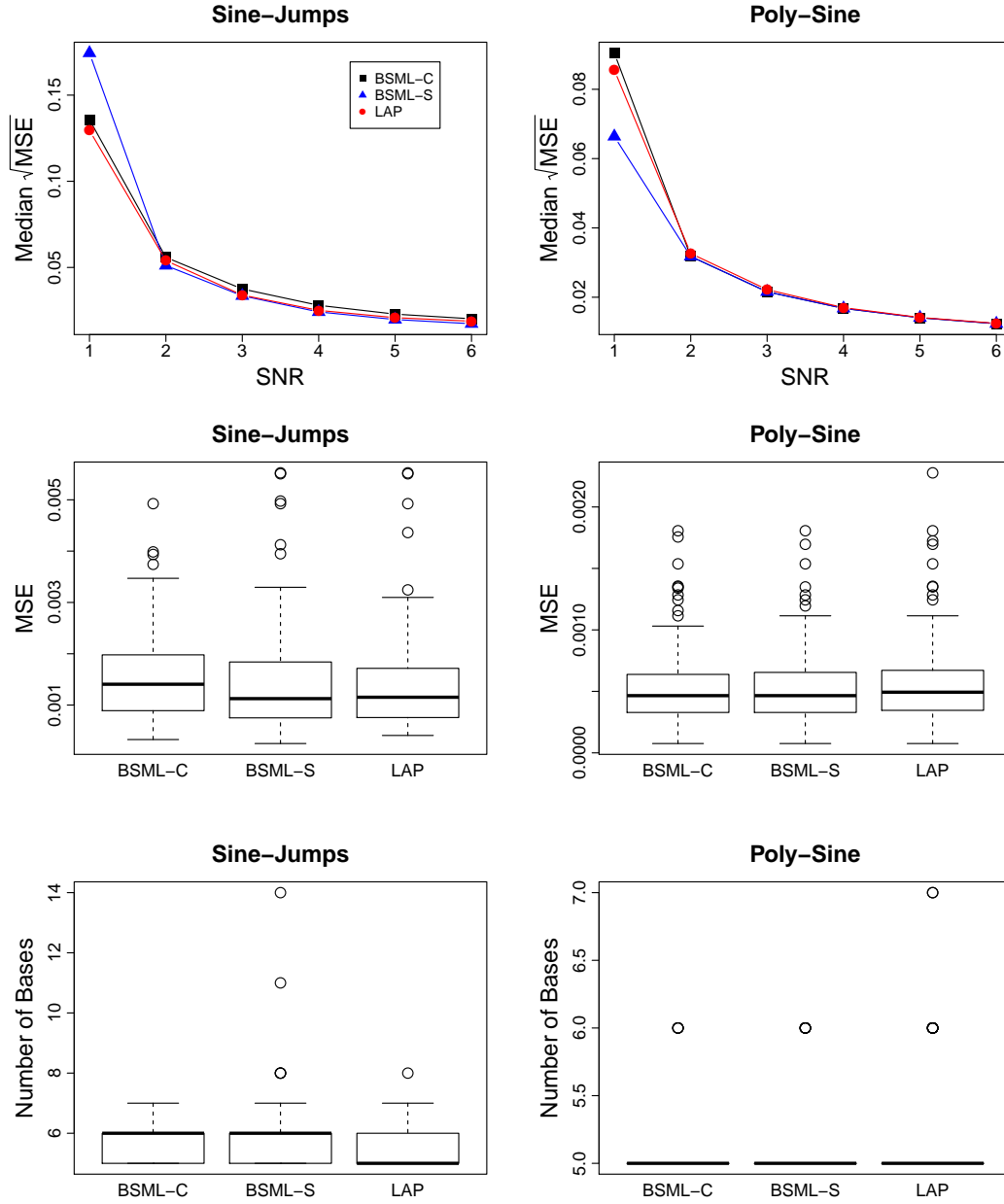


Figure 4.8: Simulation results for the Sine-Jumps and Poly-Sine examples with $n = 1024$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for both examples, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, and \bullet for the look-ahead procedure. The middle panel shows the boxplots of the MSE when SNR = 3. The bottom panel shows the number of basis functions selected when SNR = 3.

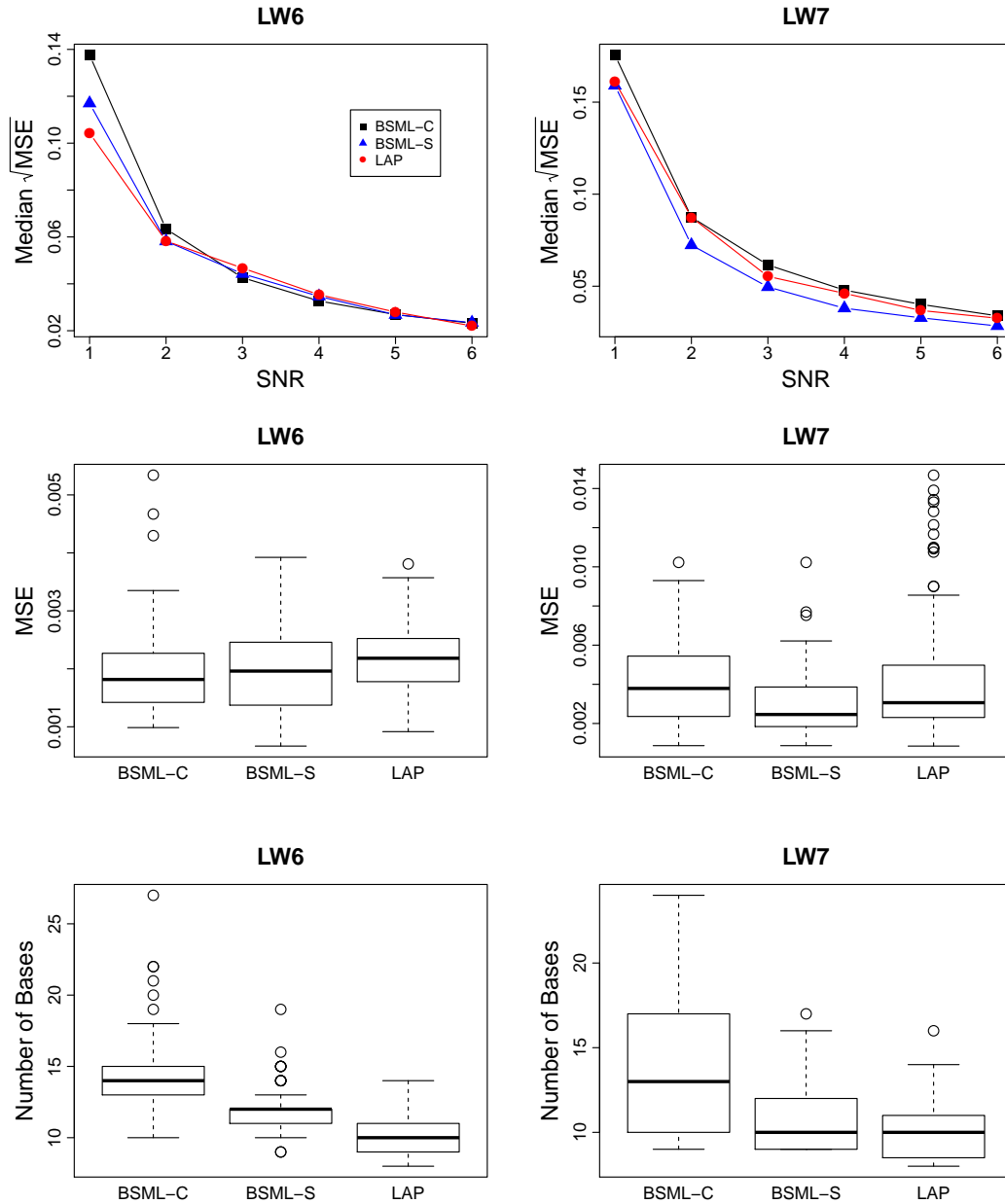


Figure 4.9: Simulation results for the LW6 and LW7 examples with $n = 1024$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for both examples, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, and \bullet for the look-ahead procedure. The middle panel shows the boxplots of the MSE when SNR = 3. The bottom panel shows the number of basis functions selected when SNR = 3.

Table 4.1: Average CPU time (in seconds) of fitting the Block-Curves and the Heavisine examples using BSML-C, BSML-S, and LAP for sample sizes $n = 256, 512, 1024$ with $\text{SNR} = 3$.

	Blocks-Curves			Heavisine		
	$n = 256$	$n = 512$	$n = 1024$	$n = 256$	$n = 512$	$n = 1024$
BSML-C	2.10	8.89	64.27	2.12	9.00	56.89
BSML-S	14.93	52.58	243.76	14.77	53.05	256.37
LAP	1.04	3.92	11.33	1.00	3.74	12.02

Table 4.2: Average CPU time (in seconds) of fitting the Sine-Jumps and the Poly-sine examples using BSML-C, BSML-S, and LAP for sample sizes $n = 256, 512, 1024$ with $\text{SNR} = 3$.

	Sine-Jumps			Poly-Sine		
	$n = 256$	$n = 512$	$n = 1024$	$n = 256$	$n = 512$	$n = 1024$
BSML-C	2.10	9.01	64.80	1.26	4.12	24.63
BSML-S	14.75	52.99	245.93	9.44	29.52	124.84
LAP	0.98	4.79	7.11	0.70	1.77	9.75

Table 4.3: Average CPU time (in seconds) of fitting the LW6 and LW7 examples using BSML-C, BSML-S, and LAP for sample sizes $n = 256, 512, 1024$ with $\text{SNR} = 3$.

	LW6			LW7		
	$n = 256$	$n = 512$	$n = 1024$	$n = 256$	$n = 512$	$n = 1024$
BSML-C	2.11	8.91	55.42	2.07	8.87	64.36
BSML-S	14.64	52.80	245.53	14.66	52.32	253.17
LAP	1.11	3.82	11.26	1.07	3.63	11.45

As we mentioned in Section 3.2, LAP only considers the pairs of basis functions of the form $\{\psi_l^-, \psi_{\cdot|l}^-\}$ for $l = 1, \dots, L_-$, where $\psi_{\cdot|l}^-$ is selected given ψ_l^- already in the model. Alternatively, the selection scheme in (3.4) can be used, where all possible pairs of unselected basis functions are considered. This alternative procedure is referred to as LAP-2. The same 100 random samples of the six univariate examples with $\text{SNR} = 3$ and $n = 512$ are used to compare the performances of LAP and LAP-2. Figure 4.10 shows the boxplots of the MSE, and Table 4.4 lists the average CPU time for both procedures.

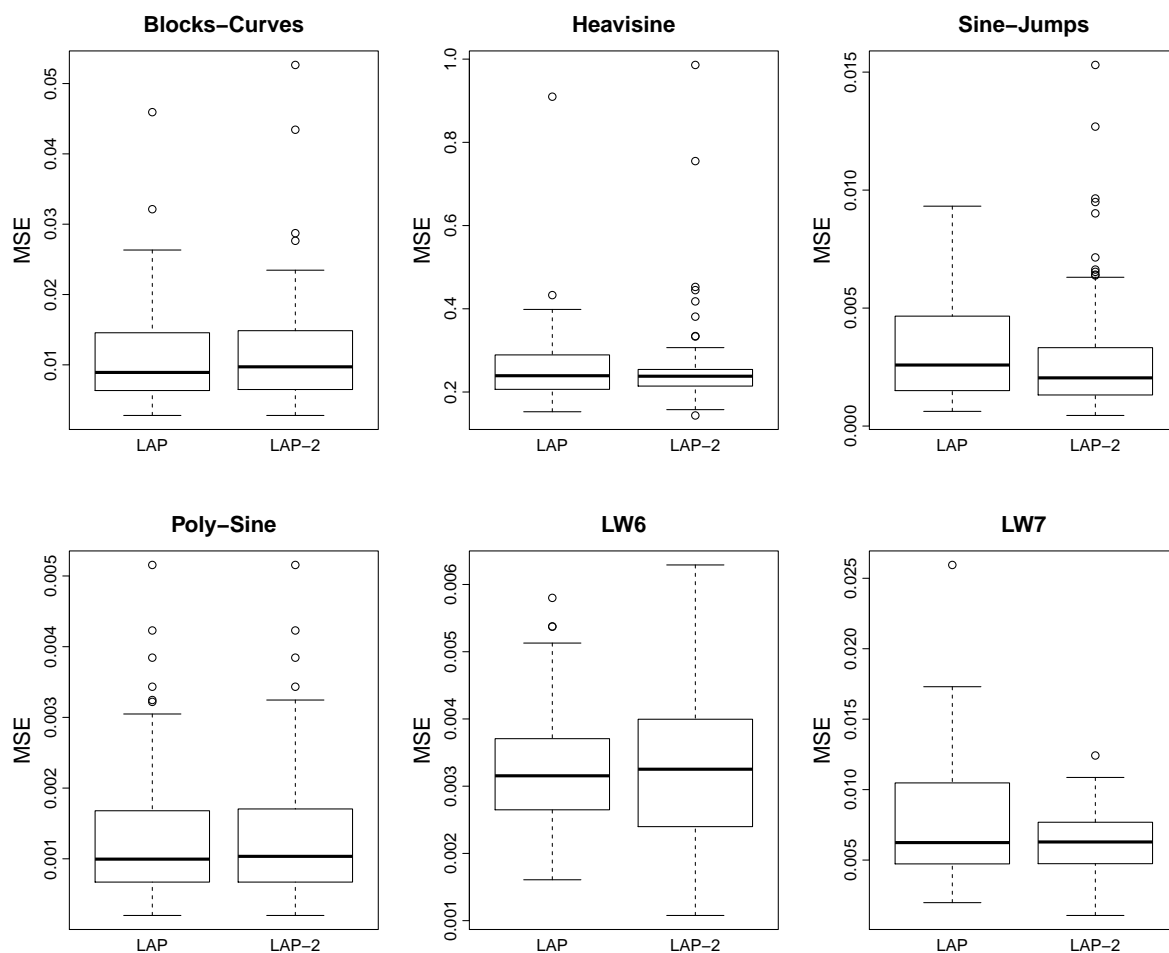


Figure 4.10: Boxplots of the MSE for fitting six univariate examples using LAP and LAP-2 with $\text{SNR} = 3$ and $n = 512$.

Table 4.4: Average CPU time (in seconds) of fitting each univariate example using LAP and LAP-2 with $n = 512$ and $\text{SNR} = 3$.

	Blocks-Curves	Heavisine	Sine-Jumps	LW6	LW7	Poly-Sine
LAP	3.92	3.74	4.79	3.82	3.63	1.77
LAP-2	124.33	64.7	142.58	127.31	109.03	13.65

We can see from Figure 4.10 and Table 4.4 that when all pairwise search scheme is used, there is no significant improvement in terms of MSE, but the computation time has increased dramatically. Therefore, we will use LAP, because it is less computationally demanding than LAP-2.

4.2 Bootstrap Confidence Intervals

We use the same Sine-Jumps and Blocks-Curves examples to evaluate the performance of the bootstrap confidence intervals from the LAP procedure, and compare them with that from the BSMML procedures. For each example, we construct the 95% bootstrap confidence intervals under two settings: $n = 256, \text{SNR} = 4$ and $n = 512, \text{SNR} = 6$ with `bootrep`=200. All the other settings are exactly the same as in 4.1. Each simulation is repeated $n_r = 100$ times.

The pointwise coverage probability (PCP) and average coverage probability (ACP) are calculated to evaluate the performance of the confidence intervals. For any n arbitrary points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ in the domain, let $f(\mathbf{x}_i)$ be the true function evaluated at \mathbf{x}_i . For the r -th simulated data set $\{(\mathbf{x}_1, y_{1,r}), \dots, (\mathbf{x}_n, y_{n,r})\}$, where $r = 1, \dots, n_r$, let $\hat{f}_{\alpha/2}^{(r)}(\mathbf{x}_i)$ and $\hat{f}_{1-\alpha/2}^{(r)}(\mathbf{x}_i)$ be the lower and upper bound of the $100(1 - \alpha)\%$ bootstrap confidence interval at \mathbf{x}_i . Then the pointwise coverage probability (PCP) at \mathbf{x}_i is defined as

$$\text{PCP}_i = \frac{1}{n_r} \sum_{r=1}^{n_r} I\left(f(\mathbf{x}_i) \in [\hat{f}_{\alpha/2}^{(r)}(\mathbf{x}_i), \hat{f}_{1-\alpha/2}^{(r)}(\mathbf{x}_i)]\right), \quad (4.1)$$

where $I(\cdot)$ is the indicator function. The average coverage probability (ACP) is computed for each $100(1 - \alpha)\%$ bootstrap confidence interval across the function, defined as

$$\text{ACP}_r = \frac{1}{n} \sum_{i=1}^n I\left(f(\mathbf{x}_i) \in [\hat{f}_{\alpha/2}^{(r)}(\mathbf{x}_i), \hat{f}_{1-\alpha/2}^{(r)}(\mathbf{x}_i)]\right). \quad (4.2)$$

Besides coverage probabilities, we are also interested in the width of a confidence interval, since given the same coverage probability a narrower confidence interval provides more precise information about the true function than a wider one. We define the pointwise average interval width (PAIW) at each design point \mathbf{x}_i as follows, where the

average is across the results for the n_r simulated data sets:

$$\text{PAIW}_i = \frac{1}{n_r} \sum_{r=1}^{n_r} \left(\hat{f}_{1-\alpha/2}^{(r)}(\mathbf{x}_i) - \hat{f}_{\alpha/2}^{(r)}(\mathbf{x}_i) \right). \quad (4.3)$$

For the Blocks-Curves and Sine-Jumps examples, Figures 4.11 to 4.14 show the comparisons of the PCP, ACP, and PAIW based on the 95% bootstrap confidence intervals from the LAP, BSML-C, and BSML-S procedures for $n_r = 100$. For all settings the BSML-S confidence intervals gave the lowest median ACP and virtually always the lowest PCP values. The median ACP for BSML-S confidence intervals for each of the four simulation cases considered is below the nominal value 0.95. The BSML-C confidence intervals and LAP confidence intervals have similar performances, and the median ACP for both of them are above the nominal value for all settings. The advantage of the LAP confidence intervals is that in general they are pointwise narrower than the BSML-C confidence intervals. For the Sine-Jumps example the LAP confidence intervals have shorter widths almost uniformly than the BSML-C confidence intervals, except at the places near 0.25 and 0.5 where the Sine-Jumps function is discontinuous.

We can also see that coverage probabilities depend on the true function, as well as the sample size and signal-to-noise ratio. Although the median ACP for all LAP confidence interval is above the nominal value of 0.95, the pointwise coverage probabilities can still be below 0.95. For the Blocks-Curves example when $n = 512$ and $\text{SNR} = 6$, the PCP is much lower than 0.95 at the design points near 0.7 where the Blocks-Curves function has a big jump. Therefore, same as the conclusions in Wang and Wahba (1995), the bootstrap confidence intervals from the LAP procedure should be interpreted as across the curve, instead of pointwise. As we mentioned earlier, the bias of the estimates tends to be large at places where f has sharp curvature or sudden jumps. Our simulations support this claim, and we found that the confidence intervals tend to be very wide at

these places as well.

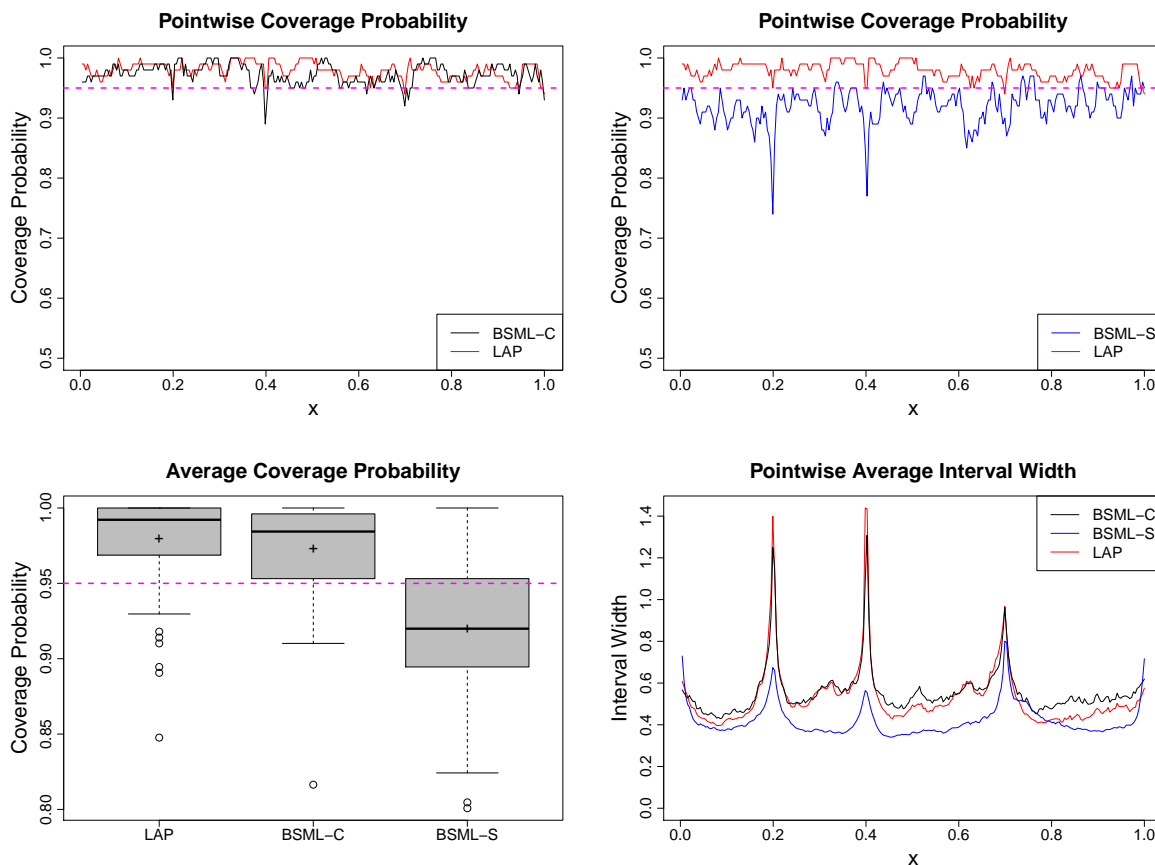


Figure 4.11: PCP, ACP, PAIW of the 95% bootstrap confidence intervals from the LAP, BSML-C, and BSML-S procedures using the Blocks-Curves example with $\text{SNR} = 4$ and $n = 256$. Top left: pointwise coverage probability comparisons between the LAP and BSML-C confidence intervals. Top right: pointwise coverage probability comparisons between the LAP and BSML-S confidence intervals. Bottom left: boxplots of the average coverage probabilities of the LAP, BSML-C, and BSML-S confidence intervals. Bottom right: pointwise average interval width of the LAP, BSML-C, and BSML-S confidence intervals. In the PCP and PAIW plots, red solid curve represents LAP, black solid curve represents BSML-C, and blue solid curve represents BSML-S. In the three coverage probability plots, the magenta dashed line represents the nominal coverage probability of 0.95. The + symbols in the ACP panel show the means for the three methods.

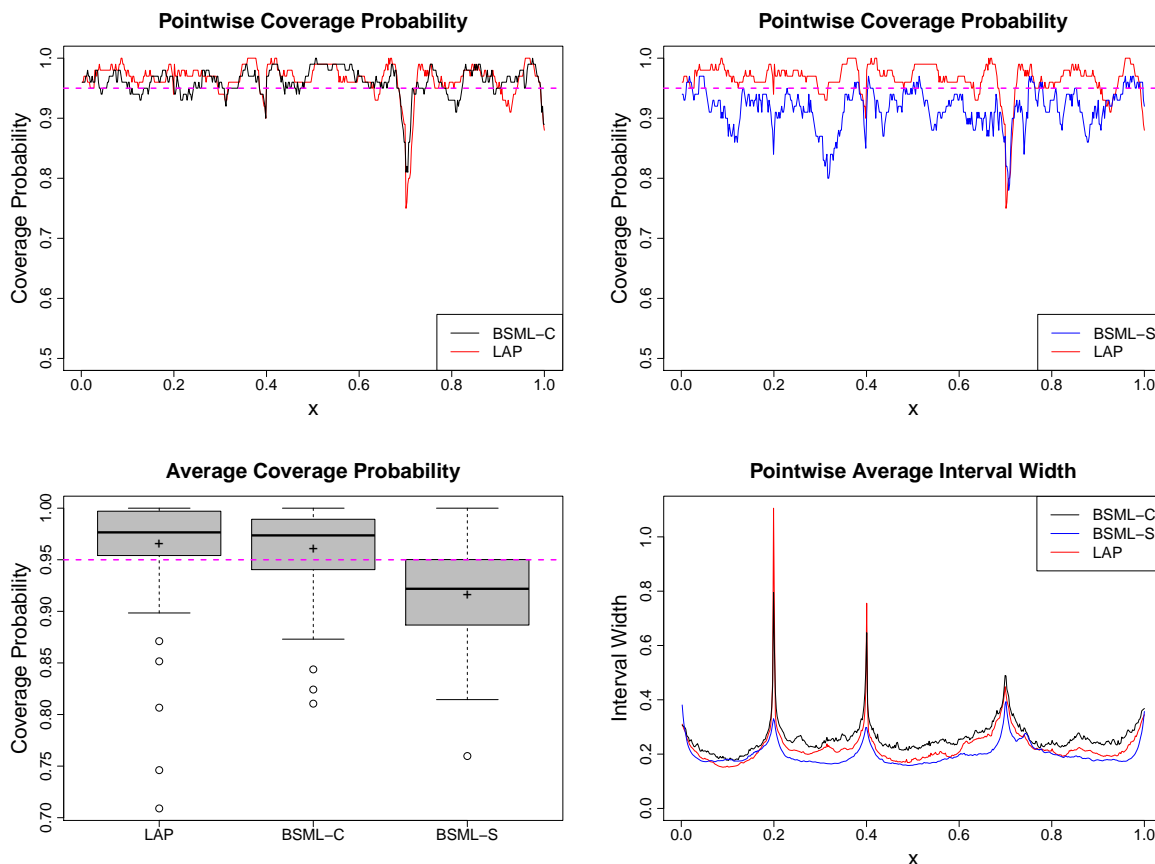


Figure 4.12: PCP, ACP, PAIW of the 95% bootstrap confidence intervals from LAP, BSML-C, and BSML-S procedures using the Blocks-Curves example with $\text{SNR} = 6$ and $n = 512$. Top left: pointwise coverage probability comparisons between the LAP and BSML-C confidence intervals. Top right: pointwise coverage probability comparisons between the LAP and BSML-S confidence intervals. Bottom left: boxplots of the average coverage probabilities of the LAP, BSML-C, and BSML-S confidence intervals. Bottom right: pointwise average interval width of the LAP, BSML-C, and BSML-S confidence intervals. In the PCP and PAIW plots, red solid curve represents LAP, black solid curve represents BSML-C, and blue solid curve represents BSML-S. In the three coverage probability plots, the magenta dashed line represents the nominal coverage probability of 0.95. The + symbols in the ACP panel show the means for the three methods.

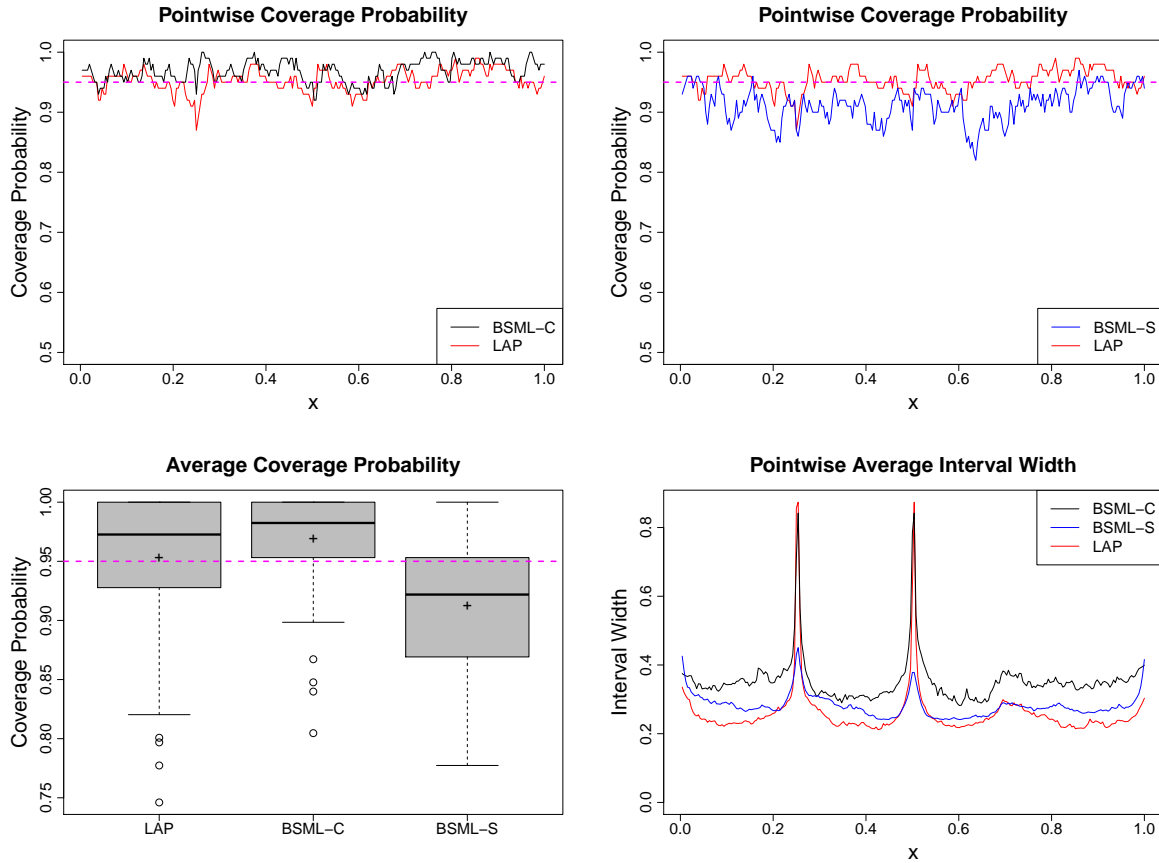


Figure 4.13: PCP, ACP, PAIW of the 95% bootstrap confidence intervals from LAP, BSML-C, and BSML-S procedures using the Sine-Jumps example with $\text{SNR} = 4$ and $n = 256$. Top left: pointwise coverage probability comparisons between the LAP and BSML-C confidence intervals. Top right: pointwise coverage probability comparisons between the LAP and BSML-S confidence intervals. Bottom left: boxplots of the average coverage probabilities of the LAP, BSML-C, and BSML-S confidence intervals. Bottom right: pointwise average interval width of the LAP, BSML-C, and BSML-S confidence intervals. In the PCP and PAIW plots, red solid curve represents LAP, black solid curve represents BSML-C, and blue solid curve represents BSML-S. In the three coverage probability plots, the magenta dashed line represents the nominal coverage probability of 0.95. The + symbols in the ACP panel show the means for the three methods.

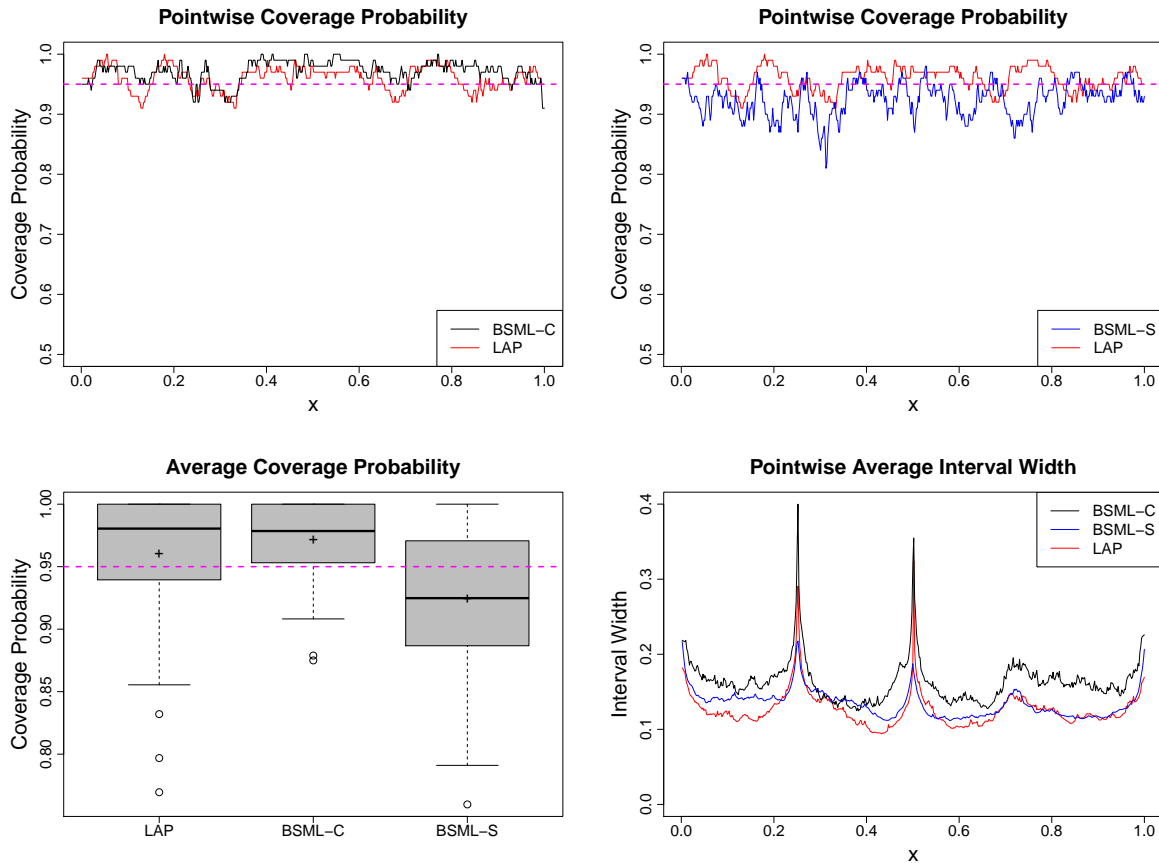


Figure 4.14: PCP, ACP, PAIW of the 95% bootstrap confidence intervals from LAP, BSML-C, and BSML-S procedures using the Sine-Jumps example with $\text{SNR} = 6$ and $n = 512$. Top left: pointwise coverage probability comparisons between the LAP and BSML-C confidence intervals. Top right: pointwise coverage probability comparisons between the LAP and BSML-S confidence intervals. Bottom left: boxplots of the average coverage probabilities of the LAP, BSML-C, and BSML-S confidence intervals. Bottom right: pointwise average interval width of the LAP, BSML-C, and BSML-S confidence intervals. In the PCP and PAIW plots, red solid curve represents LAP, black solid curve represents BSML-C, and blue solid curve represents BSML-S. In the three coverage probability plots, the magenta dashed line represents the nominal coverage probability of 0.95. The + symbols in the ACP panel show the means for the three methods.

4.3 Estimation of Multivariate Functions

4.3.1 Additive Model Examples

Now consider the bivariate nonparametric regression model

$$y_i = f_1(x_{i1}) + f_2(x_{i2}) + \varepsilon_i, \quad i = 1, \dots, n, \quad (4.4)$$

where $x_{i1} \in [0, 1]$ and $x_{i2} \in [0, 1]$ are the observations of two explanatory variables for observation y_i , and the independent random errors $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ have mean 0 and constant variance σ^2 . We choose from the same six univariate functions in Table 2.1 to construct three combinations, which are shown in Table 4.5. Table 4.6 contains the basis libraries used in BSML-C, BSML-S, and LAP for each example.

Table 4.5: Test functions for additive model simulations.

Function Name	True Function $f(\mathbf{x})$	$SD(f)$
Sine-Jumps + LW7	$\sin(2\pi x_1) - 1_{(0.5,1]}(x_1) + 1_{(0.25,1]}(x_1)$ $+ (4x_2 - 2) + 2 \exp[-256(x_2 - 0.5)^2]$	1.6141
Blocks-Curves + Sine-Jumps	$1_{[0.2,0.4)}(x_1) + \exp[(x_1 + 0.5)^2] 1_{[0.4,0.7)}(x_1) - x_1^{10} 1_{[0.7,1]}(x_1)$ $+ \sin(2\pi x_1) - 1_{(0.5,1]}(x_1) + 1_{(0.25,1]}(x_1)$	1.7702
Poly-Sine + LW6	$\sin(16\pi x_1) - 8(x_1 - 0.5)^2 + 8(x_1 - 0.5)^3 1_{(0.5,1]}(x_1)$ $+ \sin[2(4x_2 - 2)] + 2 \exp[-256(x_2 - 0.5)^2]$	1.1866

In the simulation, the design points (x_{i1}, x_{i2}) for $i = 1, \dots, n$ are uniformly sampled from $[0, 1] \times [0, 1]$, where the sample size is $n \in \{256, 512, 1024\}$. The error variance σ^2 is chosen such that $\text{SNR} \in \{1, 2, 3, 4, 5, 6\}$. The knots used in all libraries are grid points $\{0.06 + 0.02j, j = 1, 2, \dots, 44\}$. For the BSML and LAP procedures, the maximal number of bases M is fixed at 30. For LAP, the IDF spaces are $S_i = \{1, 2, 3\}$ for $i = 1, 2, 3, 4$, and M_d is set to be $\prod_{i=1}^4 |S_i| = 3^4 = 81$. The 10-fold cross validation is used to select

Table 4.6: Basis libraries used in BSML-C, BSML-S, and LAP for each additive model example. Notations of basis libraries can be found in Table 1.1.

Function Name	Basis Libraries				
	\mathcal{L}_0	\mathcal{L}_1	\mathcal{L}_2	\mathcal{L}_3	\mathcal{L}_4
Sine-Jumps + LW7	$\{1, x_1, x_2\}$	$\mathcal{P}_2(x_1)$	$\mathcal{T}_0(x_1)$	$\mathcal{C}_2(x_2)$	$\mathcal{T}_2(x_2)$
Blocks-Curves + Sine-Jumps	$\{1, x_1, x_2\}$	$\mathcal{C}_2(x_1)$	$\mathcal{T}_0(x_1)$	$\mathcal{P}_2(x_2)$	$\mathcal{T}_0(x_2)$
Poly-Sine + LW6	$\{1, x_1, x_2\}$	$\mathcal{F}_{8,25}(x_1)$	$\mathcal{T}_2(x_1)$	$\mathcal{C}_2(x_2)$	$\mathcal{T}_2(x_2)$

the final model in the look-ahead procedure.

We also applied the SS-ANOVA model, the COSSO method, and the MARS procedure to these three examples. We used the function `ssr` in the R package `assist` (Wang and Ke (2015)) to fit the SS-ANOVA models. For the Blocks-Curves + Sine-Jumps example, we assume $f \in W_2^2[0, 1] \otimes W_2^2(per)$. For the Sine-Jumps + LW7 example, we assume $f \in W_2^2(per) \otimes W_2^2[0, 1]$. For the Poly-Sine + LW6 example, we assume $f \in W_2^2[0, 1] \otimes W_2^2[0, 1]$. For all SS-ANOVA models, only the main effects are included. For the COSSO method, we used the function `cosso` in the R package `cosso` (Zhang and Lin (2013)), with 10-fold cross validation for the tuning parameter selection. For the MARS procedure, we used the function `mars` in the R package `mda` (Hastie et al. (2015)), with `degree = 1` which means no interaction term is included. The IDF used in `mars` is fixed at 3, and the maximum number of model terms `nk` is fixed at 30. The simulation results for BSML-C, BSML-S, LAP, SS-ANOVA, COSSO, and MARS based on 100 random samples are shown in Figures 4.15 – 4.17. The average CPU time for each sample size is listed in Table 4.7.

Based on the results, MARS did not work well when $SNR \geq 2$ for all three examples. For COSSO, it did not work well for the Poly-Sine + LW6 example. When $SNR \geq 4$, LAP has the best performance for the Blocks-Curves + Sine-Jumps and Polysine + LW6 examples across all sample sizes, and also for the Sine-Jumps + LW7 example when

the sample size is 256 or 512. On average LAP has selected fewer basis functions than BSML-C and BSML-S. SS-ANOVA has comparable performance as the BSML and LAP procedures, but in terms of computation time it has disadvantage when the sample size is big, e.g., when $n = 1024$. LAP and BSML-C have similar computation speed, both are much faster than BSML-S.

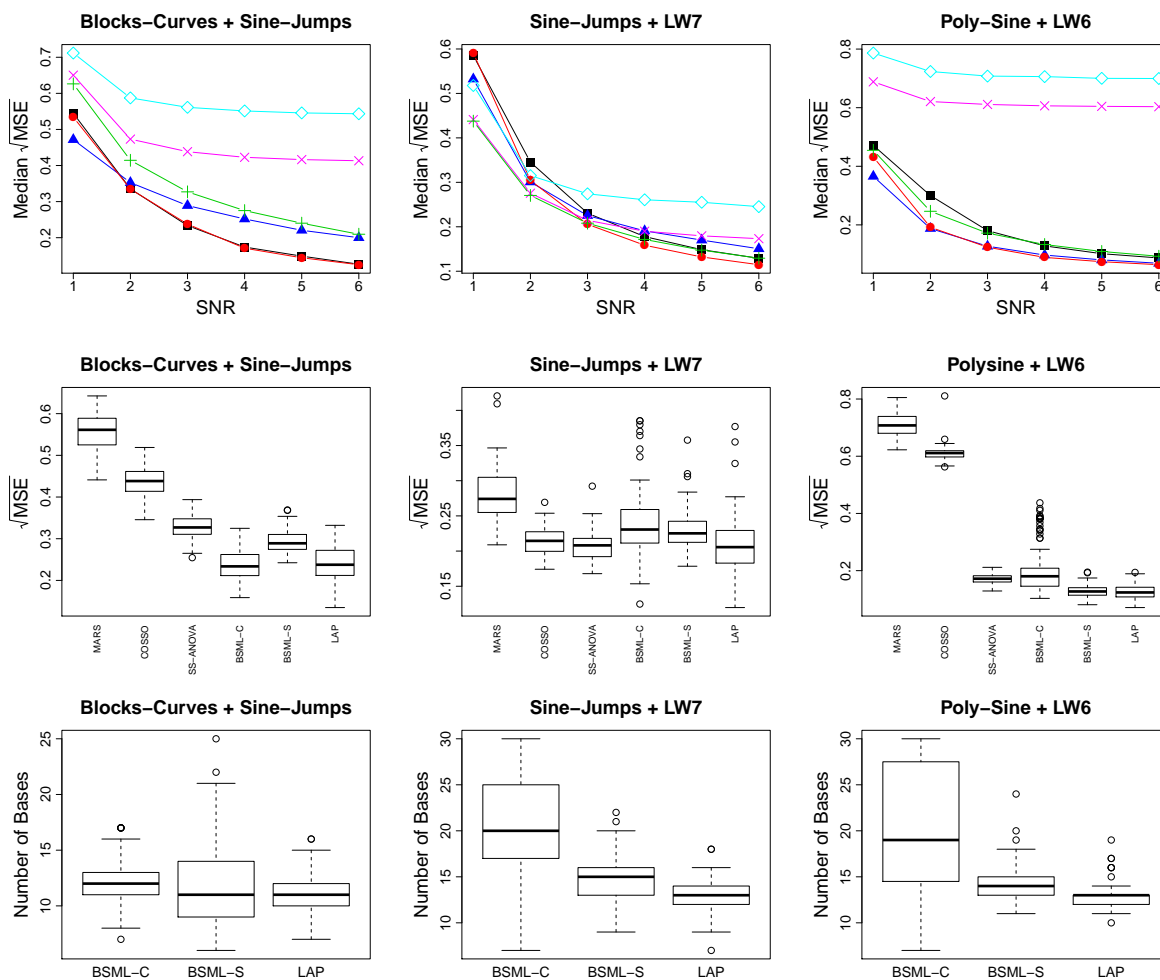


Figure 4.15: Simulation results for three additive model examples with $n = 256$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for each example, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, \bullet for LAP, $+$ for SS-ANOVA, \times for COSO, and \diamond for MARS. The middle panel shows the boxplots of the $\sqrt{\text{MSE}}$ when SNR = 3. The bottom panel shows the number of basis functions selected by BSML-C, BSML-S, and LAP when SNR = 3.

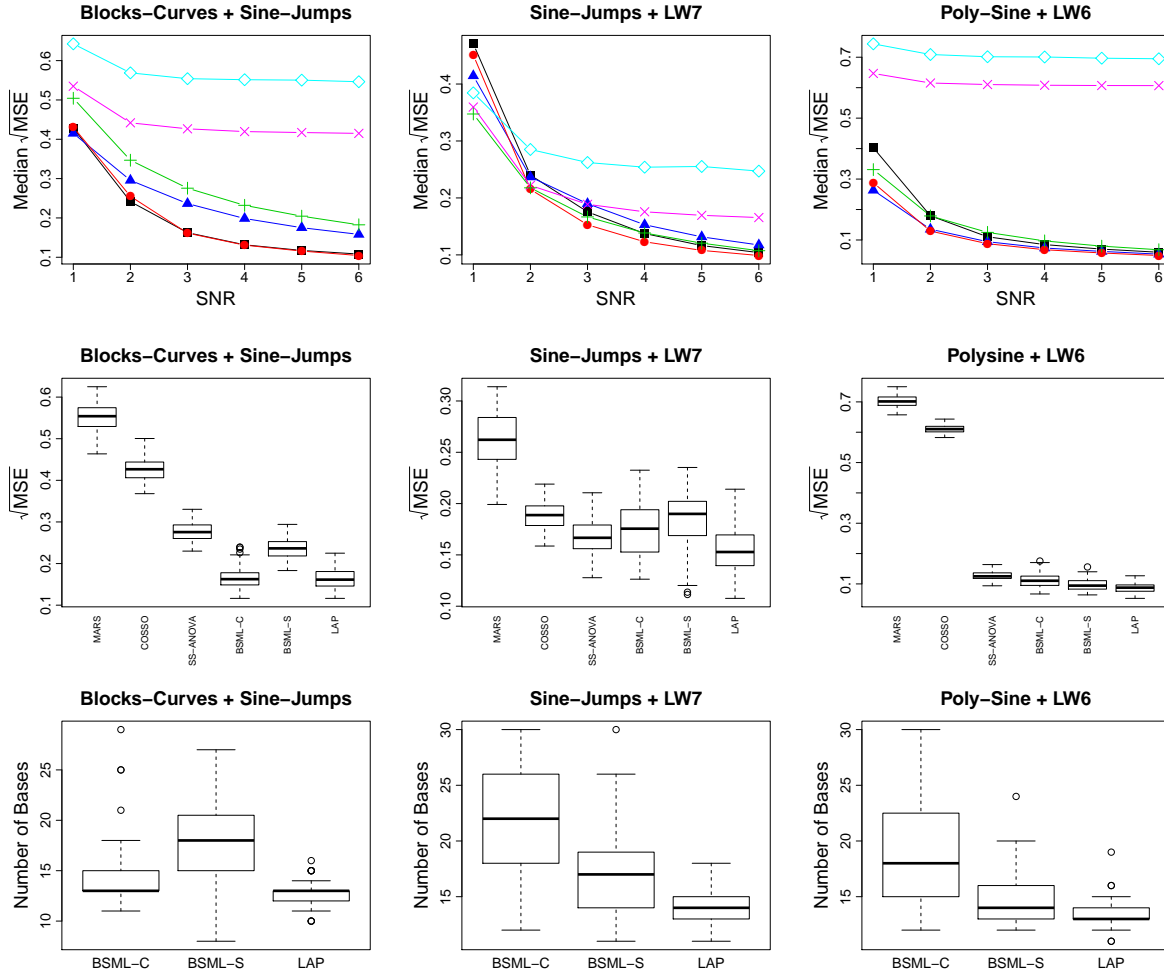


Figure 4.16: Simulation results for three additive model examples with $n = 512$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for each example, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, \bullet for LAP, $+$ for SS-ANOVA, \times for COSSO, and \diamond for MARS. The middle panel shows the boxplots of the $\sqrt{\text{MSE}}$ when $\text{SNR} = 3$. The bottom panel shows the number of basis functions selected by BSML-C, BSML-S, and LAP when $\text{SNR} = 3$.

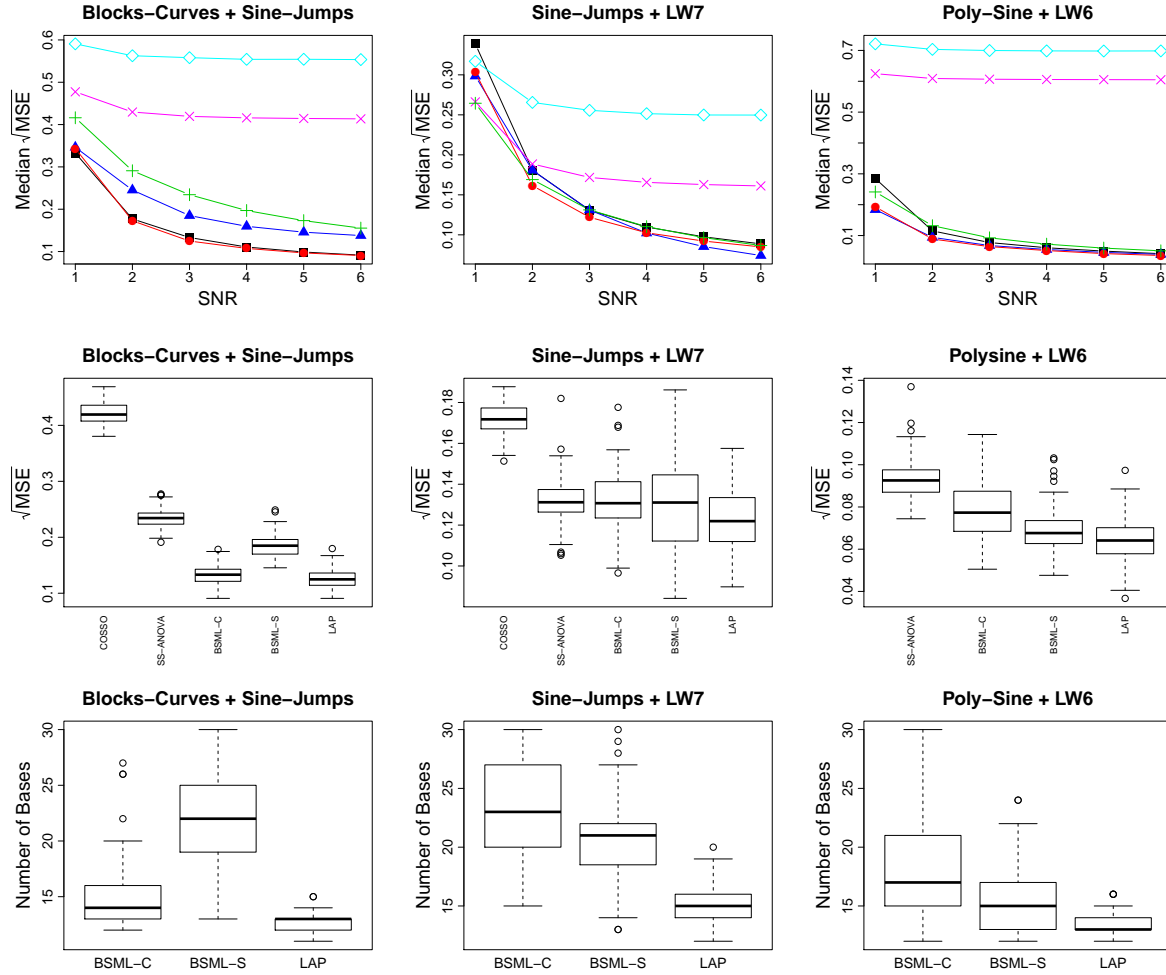


Figure 4.17: Simulation results for three additive model examples with $n = 1024$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for each example, where the symbols used are: ■ for BSML-C, ▲ for BSML-S, ● for LAP, + for SS-ANOVA, × for COSSO, and ◇ for MARS. The middle panel shows the boxplots of the $\sqrt{\text{MSE}}$ when SNR = 3. The results of fitting all three examples with MARS and the result of fitting the Polysine + LW6 example with COSSO are not shown here because they are much worse compared to the results from the BSML and LAP procedures. The bottom panel shows the number of basis functions selected by BSML-C, BSML-S, and LAP when SNR = 3.

Table 4.7: Average CPU time (in seconds) of fitting each additive model example using BSML-C, BSML-S, LAP, SS-ANOVA, COSSO, MARS for sample sizes $n = 256, 512, 1024$ with $\text{SNR} = 3$.

	Blocks-Curves + Sine-Jumps			Sine-Jumps + LW7		
	$n = 256$	$n = 512$	$n = 1024$	$n = 256$	$n = 512$	$n = 1024$
BSML-C	0.89	1.73	3.51	1.04	1.35	3.29
BSML-S	8.23	16.26	30.82	8.16	22.23	37.46
LAP	0.65	1.37	2.04	1.09	1.86	2.85
SS-ANOVA	0.32	1.66	9.42	0.18	1.00	6.87
COSSO	0.39	0.87	3.11	0.27	0.86	2.61
MARS	0.05	0.06	0.05	0.02	0.02	0.02

	Poly-Sine + LW6		
	$n = 256$	$n = 512$	$n = 1024$
BSML-C	1.04	1.39	3.33
BSML-S	8.36	22.80	39.34
LAP	0.81	1.18	1.56
SS-ANOVA	0.20	1.05	6.00
COSSO	0.30	0.94	2.88
MARS	0.02	0.02	0.03

4.3.2 Models with Interactions

Consider the bivariate nonparametric regression model

$$y_i = f_{12}(x_{i1}, x_{i2}) + \varepsilon_i, \quad i = 1, \dots, n, \quad (4.5)$$

where $x_{i1} \in [0, 1]$ and $x_{i2} \in [0, 1]$ are the observations of two explanatory variables for observation y_i , and the independent random errors $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ have mean 0 and constant variance σ^2 . We constructed three examples of such function f , which are shown in Table 4.8. Table 4.9 and 4.10 contain the basis libraries used in the BSML and LAP procedures for each example. Here the knots used are $z_j = 0.06 + 0.02j$, $j = 1, \dots, 44$. For all three examples, \mathcal{L}_1 to \mathcal{L}_5 contain similar spline representers used in SS-ANOVA model with the model space $W_2^m[0, 1] \otimes W_2^m[0, 1]$. \mathcal{L}_6 for the SJ·LW7 and BC·LW6 examples contains 2d step functions in x_1 , which are used to capture the jumps in the true functions.

Table 4.8: Test functions for models with interactions.

Function Name	True Function $f(\mathbf{x})$	$SD(f)$
SJ·LW7	$2[0.5 \sin(2\pi x_1) - 1_{(0.5,1]}(x_1) + 1_{(0.24,1]}(x_1)]$ $+4x_2 + 6(x_1 - 0.5) \exp[-256(x_2 - 0.5)^2]$	1.7930
BC·LW6	$1_{[0.2,0.4)}(x_1) + \exp[(x_1 + 0.5)^2]1_{[0.4,0.7)}(x_1) - x_1^{10}1_{[0.7,1]}(x_1)$ $+ \sin[2(4x_2 - 2)] + 6(x_1 - 0.5) \exp[-256(x_2 - 0.5)^2]$	1.6671
LW6·LW7	$\sin[2(4x_1 - 2)] + 6(x_2 - 0.5) \exp[-256(x_1 - 0.5)^2]$ $+4(x_2 - 0.5) + 2 \exp[-256(x_2 - 0.5)^2]$	1.6483

In the simulation, the design points (x_{i1}, x_{i2}) for $i = 1, \dots, n$ are uniformly sampled from $[0, 1] \times [0, 1]$ for all examples. The sample size used is $n \in \{256, 512, 1024\}$. The error variance σ^2 is chosen such that $\text{SNR} \in \{1, 2, 3, 4, 5, 6\}$. For the BSML and LAP procedures, the maximal number of bases M is fixed at 40. For the look-ahead procedure,

the IDF spaces are $S_i = \{1, 2, 3\}$ for each library \mathcal{L}_i . M_d is set to be $\prod_{i=1}^6 |S_i| = 3^6 = 729$ for the SJ·LW7 and BC·LW6 examples, and $\prod_{i=1}^5 |S_i| = 3^5 = 243$ for the LW6·LW7 example. The 10-fold cross validation is used to select the final model in the look-ahead procedure.

Table 4.9: Basis libraries for the SJ·LW7 and BC·LW6 examples.

Basis libraries for the SJ·LW7 and BC·LW6 examples
$\mathcal{L}_0 = \{1, x_1 - 0.5, x_2 - 0.5, (x_1 - 0.5)(x_2 - 0.5)\}$
$\mathcal{L}_1 = \mathcal{C}_2(x_1) = \left\{ \{R_2(x_{i1}, z_1)\}_{i=1}^n, \dots, \{R_2(x_{i1}, z_q)\}_{i=1}^n \right\}$
$\mathcal{L}_2 = \mathcal{C}_2(x_2) = \left\{ \{R_2(x_{i2}, z_1)\}_{i=1}^n, \dots, \{R_2(x_{i2}, z_q)\}_{i=1}^n \right\}$
$\mathcal{L}_3 = \left\{ \{(z_1 - 0.5)(x_{i1} - 0.5)R_2(x_{i2}, z_1)\}_{i=1}^n, \dots, \{(z_q - 0.5)(x_{i1} - 0.5)R_2(x_{i2}, z_q)\}_{i=1}^n \right\}$
$\mathcal{L}_4 = \left\{ \{(z_2 - 0.5)(x_{i2} - 0.5)R_2(x_{i1}, z_1)\}_{i=1}^n, \dots, \{(z_q - 0.5)(x_{i2} - 0.5)R_2(x_{i1}, z_q)\}_{i=1}^n \right\}$
$\mathcal{L}_5 = \left\{ \{R_2(x_{i1}, z_1)R_2(x_{i2}, z_1)\}_{i=1}^n, \dots, \{R_2(x_{i1}, z_q)R_2(x_{i2}, z_q)\}_{i=1}^n \right\}$
$\mathcal{L}_6 = \mathcal{T}_0(x_1) = \left\{ \{1_{(z_1, 1]}(x_{i1})\}_{i=1}^n, \dots, \{1_{(z_q, 1]}(x_{i1})\}_{i=1}^n \right\}$

Table 4.10: Basis libraries for the LW6·LW7 example.

Basis libraries for the LW6·LW7 example
$\mathcal{L}_0 = \{1, x_1 - 0.5, x_2 - 0.5, (x_1 - 0.5)(x_2 - 0.5)\}$
$\mathcal{L}_1 = \mathcal{C}_2(x_1) = \left\{ \{R_2(x_{i1}, z_1)\}_{i=1}^n, \dots, \{R_2(x_{i1}, z_q)\}_{i=1}^n \right\}$
$\mathcal{L}_2 = \mathcal{C}_2(x_2) = \left\{ \{R_2(x_{i2}, z_1)\}_{i=1}^n, \dots, \{R_2(x_{i2}, z_q)\}_{i=1}^n \right\}$
$\mathcal{L}_3 = \left\{ \{(z_1 - 0.5)(x_{i1} - 0.5)R_2(x_{i2}, z_1)\}_{i=1}^n, \dots, \{(z_q - 0.5)(x_{i1} - 0.5)R_2(x_{i2}, z_q)\}_{i=1}^n \right\}$
$\mathcal{L}_4 = \left\{ \{(z_2 - 0.5)(x_{i2} - 0.5)R_2(x_{i1}, z_1)\}_{i=1}^n, \dots, \{(z_q - 0.5)(x_{i2} - 0.5)R_2(x_{i1}, z_q)\}_{i=1}^n \right\}$
$\mathcal{L}_5 = \left\{ \{R_2(x_{i1}, z_1)R_2(x_{i2}, z_1)\}_{i=1}^n, \dots, \{R_2(x_{i1}, z_q)R_2(x_{i2}, z_q)\}_{i=1}^n \right\}$

We also applied the SS-ANOVA model and the MARS procedure to these three examples. For all examples, SS-ANOVA model assumes $f \in W_2^2[0, 1] \otimes W_2^2[0, 1]$. For the MARS procedure, we set `degree = 2` in `mars` which means we allow pairwise interactions terms in the model. The IDF used is fixed at 3, and the maximum number of model terms `nk` is fixed at 40. The simulation results for BSML-C, BSML-S, LAP, SS-ANOVA, and MARS based on 100 random samples are shown in Figure 4.18 – 4.20. The average CPU time for each sample size is listed in Table 4.11.

Based on the results, for both the SJ·LW7 and BC·LW6 examples LAP has the best performance across all sample sizes when $\text{SNR} \geq 3$. For these two examples, it makes sense that the LAP and BSML procedures outperform SS-ANOVA, since they can select the step functions from \mathcal{L}_6 to fit the jumps perfectly while SS-ANOVA can only smooth out the jumps using the cubic spline representers. For the LW6·LW7 example, however, because the true function is continuous and can be approximated well by using only the cubic spline representers, SS-ANOVA has better performance than the LAP and BSML procedures. For all three examples, on average LAP selected fewer basis functions than BSML-C and BSML-S. In terms of the computation speed, LAP is not as fast as BSML-C, but is still at least 3 times faster than BSML-S.

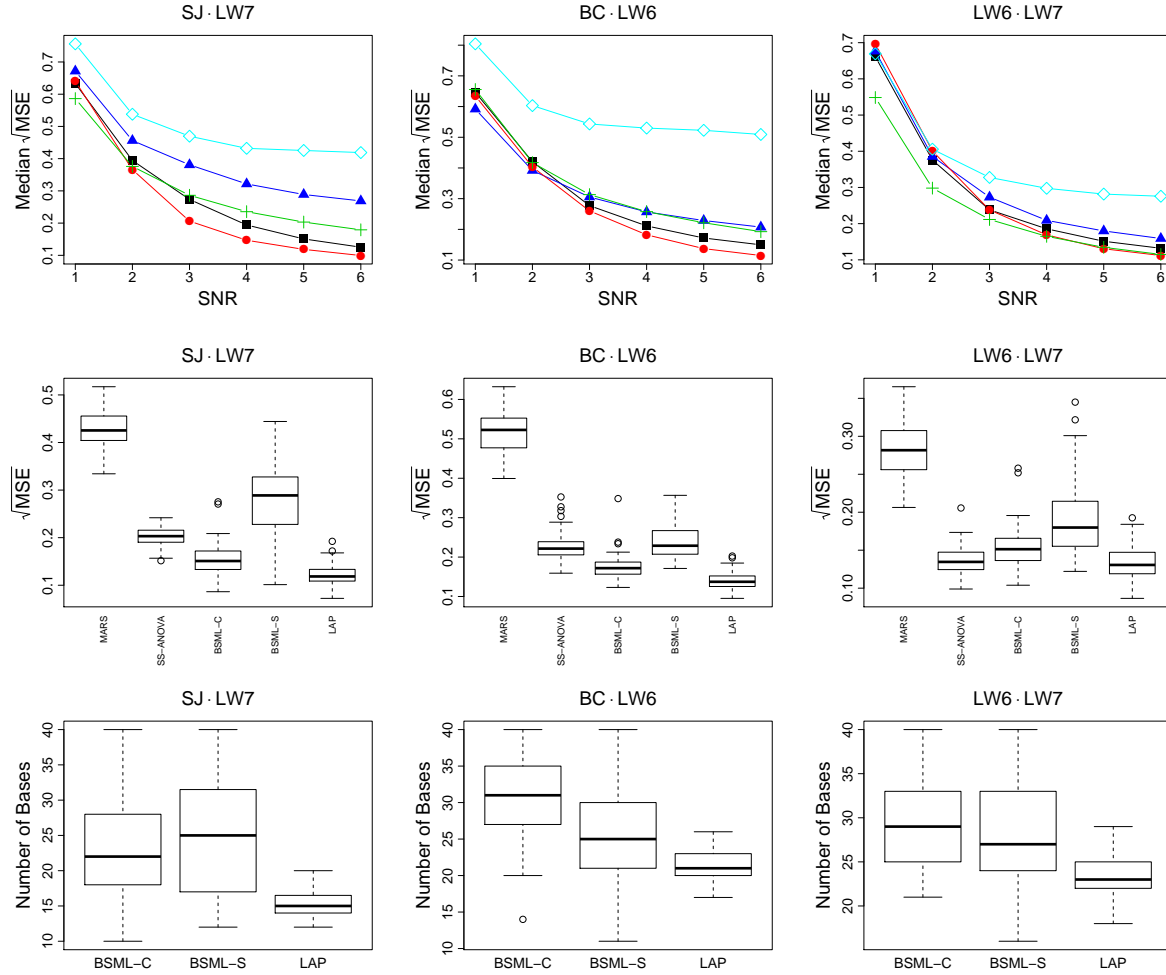


Figure 4.18: Simulation results for the SJ-LW7, BC-LW6, and LW6-LW7 examples with $n = 256$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for each example, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, \bullet for LAP, $+$ for SS-ANOVA, and \diamond for MARS. The middle panel shows the boxplots of the $\sqrt{\text{MSE}}$ when SNR = 5. The bottom panel shows the number of basis functions selected by BSML-C, BSML-S, and LAP when SNR = 5.

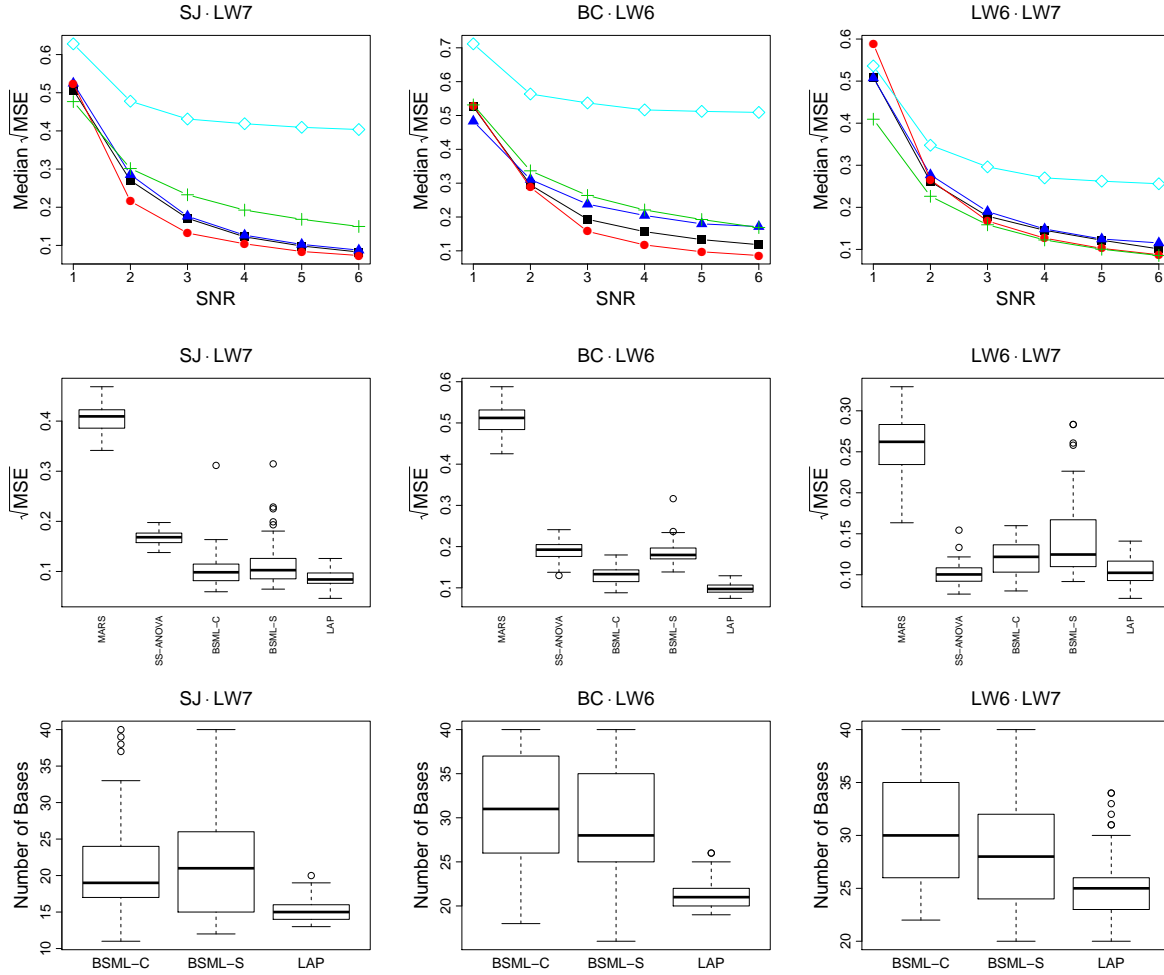


Figure 4.19: Simulation results for the SJ-LW7, BC-LW6, and LW6-LW7 examples with $n = 512$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for each example, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, \bullet for LAP, $+$ for SS-ANOVA, and \diamond for MARS. The middle panel shows the boxplots of the $\sqrt{\text{MSE}}$ when SNR = 5. The bottom panel shows the number of basis functions selected by BSML-C, BSML-S, and LAP when SNR = 5.

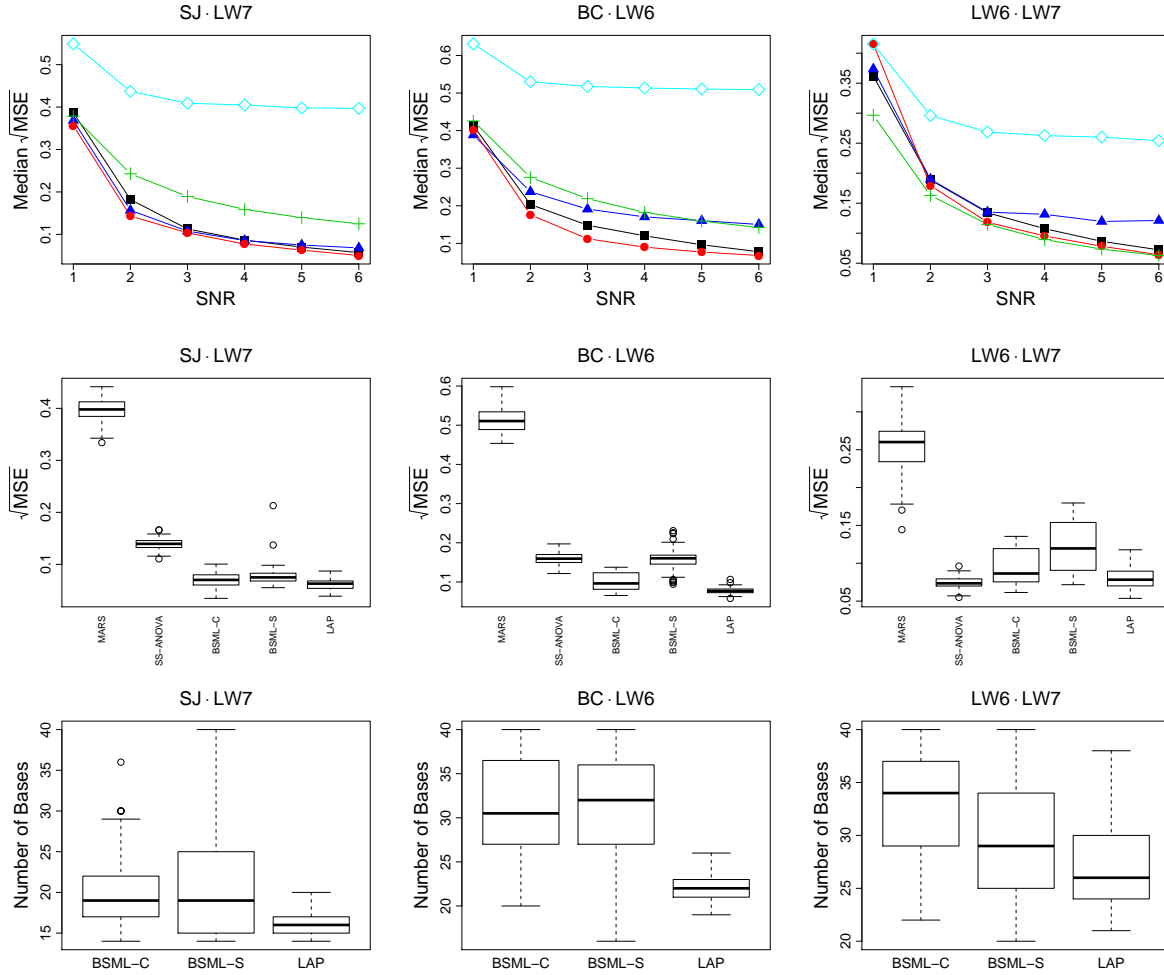


Figure 4.20: Simulation results for the SJ-LW7, BC-LW6, and LW6-LW7 examples with $n = 1024$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for each example, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, \bullet for LAP, $+$ for SS-ANOVA, and \diamond for MARS. The middle panel shows the boxplots of the $\sqrt{\text{MSE}}$ when SNR = 5. The bottom panel shows the number of basis functions selected by BSML-C, BSML-S, and LAP when SNR = 5.

Table 4.11: Average CPU time (in seconds) of fitting the SJ-LW7, BC-LW6, and LW6-LW7 examples using BSML-C, BSML-S, LAP, SS-ANOVA, MARS for sample sizes $n = 256, 512, 1024$ with $\text{SNR} = 5$.

	SJ-LW7			BC-LW6		
	$n = 256$	$n = 512$	$n = 1024$	$n = 256$	$n = 512$	$n = 1024$
BSML-C	1.67	3.10	5.18	1.61	3.19	5.20
BSML-S	17.90	34.25	67.16	17.62	34.11	67.55
LAP	3.43	5.00	8.06	4.34	7.73	11.43
SS-ANOVA	1.60	9.65	61.28	1.62	9.59	61.56
MARS	0.02	0.04	0.06	0.02	0.03	0.05

	LW6-LW7		
	$n = 256$	$n = 512$	$n = 1024$
BSML-C	1.32	2.70	4.31
BSML-S	14.74	28.39	56.02
LAP	1.74	3.34	5.46
SS-ANOVA	1.37	7.84	51.29
MARS	0.02	0.03	0.05

4.4 Three Other Examples

Now we consider three other examples where the true function f is not additive in nature. Similar examples can be found in Qiu (2005). Again consider the model in (4.5). The true functions $f(x_1, x_2)$ are listed in Table 4.12 and plotted in Figure 4.21 – 4.23. Table 4.13 – 4.15 contain the basis libraries used in the BSMML and LAP procedures for each example.

Table 4.12: Test functions of the Paraboloid-Jump, Mexican hat, and Sinusoid-Jump examples.

Function Name	True Function $f(\mathbf{x})$	$SD(f)$
Paraboloid-Jump	$-2[(x_1 - 0.5)^2 + (x_2 - 0.5)^2]$ $+0.5 \times 1_{[0,0.09)}\{(x_1 - 0.5)^2 + (x_2 - 0.5)^2\}$	0.4051
Mexican hat	$(2\pi)^{-1}(2 - x_1^2 - x_2^2) \exp(-0.5(x_1^2 + x_2^2))$	0.0498
Sinusoid-Jump	$0.5(1 - x_1)x_2 + [1 + 0.2 \sin(2\pi x_1)] 1_{[0.6 \sin(\pi x_1) + 0.2, 1]}(x_2)$	0.5716

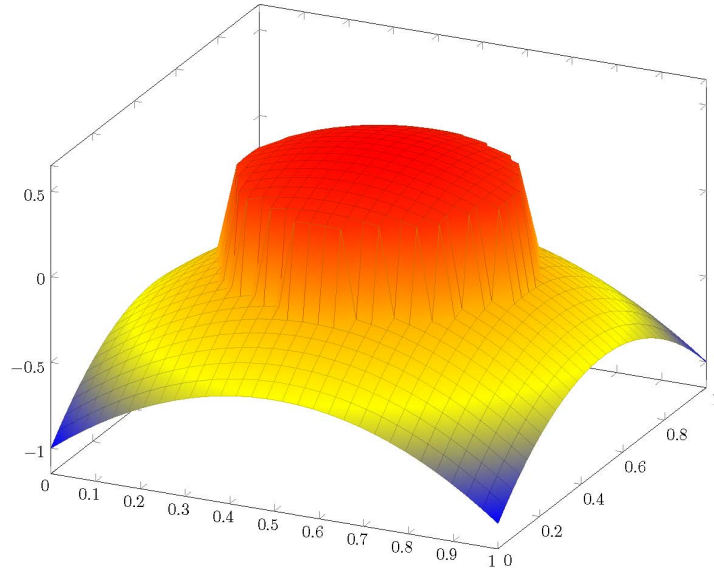


Figure 4.21: Paraboloid-Jump function.

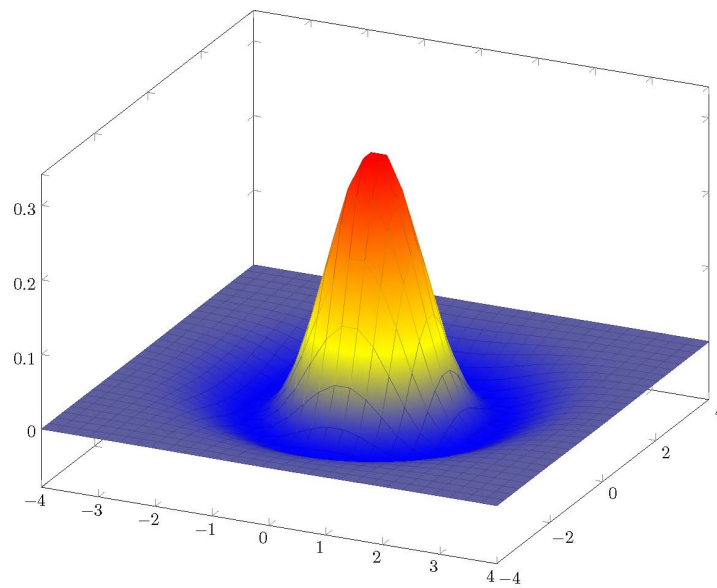


Figure 4.22: Mexican hat function.

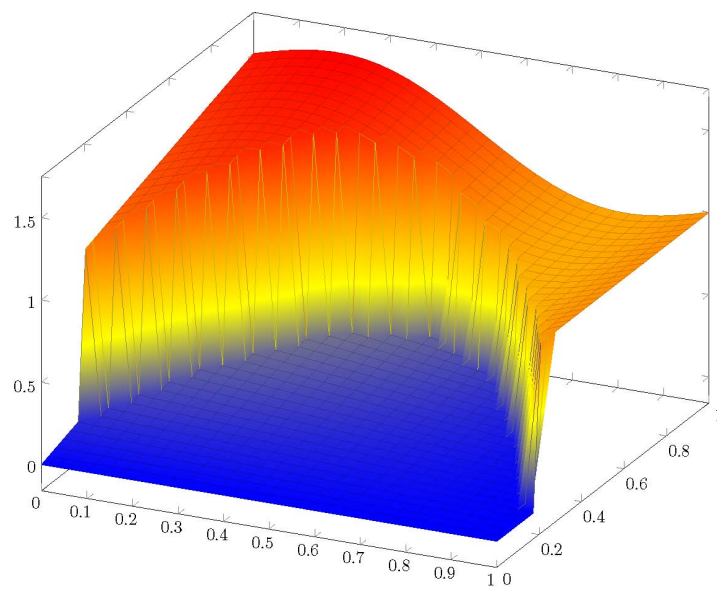


Figure 4.23: Sinusoid-Jump function.

Table 4.13: Basis libraries for the Paraboloid-Jump example.

Basis libraries for the Paraboloid-Jump example
$\mathcal{L}_0 = \{1\}$
$\mathcal{L}_1 = \{x_1, x_1^2, x_1^3\}$
$\mathcal{L}_2 = \{x_2, x_2^2, x_2^3\}$
$\mathcal{L}_3 = \{1_{[0,r]} \{(x_1 - c_1)^2 + (x_2 - c_2)^2\}, \text{ where } c_1, c_2 \in \{0.2, 0.3, \dots, 0.8\},$ $r \in \{0.1, 0.2, \dots, 0.5\}\}$
$\mathcal{L}_4 = \left\{ \exp \left\{ -0.09[(x_1 - c_1)^2 + (x_2 - c_2)^2] \right\}, \text{ where } c_1, c_2 \in \{0.2, 0.3, \dots, 0.8\} \right\}$
$\mathcal{L}_5 = \left\{ \exp \left\{ -0.36[(x_1 - c_1)^2 + (x_2 - c_2)^2] \right\}, \text{ where } c_1, c_2 \in \{0.2, 0.3, \dots, 0.8\} \right\}$

Table 4.14: Basis libraries for the Mexican hat example.

Basis libraries for the Mexican hat example
$\mathcal{L}_0 = \{1\}$
$\mathcal{L}_1 = \left\{ \exp \left\{ -0.09[(x_1 - c_1)^2 + (x_2 - c_2)^2] \right\}, \text{ where } c_1, c_2 \in \{-3, -2, -1, 0, 1, 2, 3\} \right\}$
$\mathcal{L}_2 = \left\{ \exp \left\{ -0.36[(x_1 - c_1)^2 + (x_2 - c_2)^2] \right\}, \text{ where } c_1, c_2 \in \{-3, -2, -1, 0, 1, 2, 3\} \right\}$
$\mathcal{L}_3 = \left\{ \exp \left\{ -[(x_1 - c_1)^2 + (x_2 - c_2)^2] \right\}, \text{ where } c_1, c_2 \in \{-3, -2, -1, 0, 1, 2, 3\} \right\}$
$\mathcal{L}_4 = \left\{ \exp \left\{ -4[(x_1 - c_1)^2 + (x_2 - c_2)^2] \right\}, \text{ where } c_1, c_2 \in \{-3, -2, -1, 0, 1, 2, 3\} \right\}$
$\mathcal{L}_5 = \left\{ \exp \left\{ -9[(x_1 - c_1)^2 + (x_2 - c_2)^2] \right\}, \text{ where } c_1, c_2 \in \{-3, -2, -1, 0, 1, 2, 3\} \right\}$

Table 4.15: Basis libraries for the Sinusoid-Jump example.

Basis libraries for the Sinusoid-Jump example
$\mathcal{L}_0 = \{1\}$
$\mathcal{L}_1 = \{x_1, x_2, x_1x_2\}$
$\mathcal{L}_2 = \{\sin(\pi kx_1), \cos(\pi kx_1) : k = 1, 2, 3, 4, 5\}$
$\mathcal{L}_3 = 1_{[a \sin(2\pi b^{-1}x_1)+c, 1]}(x_2)$ where $a, c \in \{0.2, 0.4, 0.6\}, b \in \{1, 2, 4\}$
$\mathcal{L}_4 = \{\xi_\ell \xi_\kappa, \text{ for } \xi_\ell \in \mathcal{L}_2, \xi_\kappa \in \mathcal{L}_3\}$ (tensor product basis functions)
$\mathcal{L}_5 = \{\xi_i \xi_j, \text{ where } \xi_i \in \{1_{[d_i, d_{i+1}]}(x_1), i = 1, \dots, 10\}, \xi_j \in \{1_{[d_j, d_{j+1}]}(x_2), j = 1, \dots, 10\},$ $d_i, d_j \text{ are the } i\text{th and } j\text{th elements of the sequence } (0, 0.1, \dots, 0.9, 1)\}$
$\mathcal{L}_6 = \mathcal{C}_2(x_1)$, with knots $\{0.02 + 0.01j, j = 1, \dots, 95\}$

For the Paraboloid-Jump example, \mathcal{L}_3 contains the 2d step functions which have jumps for points inside a circle with different center locations and radii. Bases from \mathcal{L}_3 are used to capture the jump in the Paraboloid-Jump function. \mathcal{L}_4 and \mathcal{L}_5 contains radial basis functions which have shapes close to the paraboloid.

For the Mexican hat example, the nonnull libraries contains radial basis functions with different centers and shape parameters. These bases can be used to approximate the Mexican hat function well because both are closely related to the Gaussian density.

For the Sinusoid-Jump example, \mathcal{L}_3 contains the 2d step functions which have jumps for x_2 is greater or equal to the sine functions of x_1 with different amplitudes, periods, and shifts. These basis functions are used to capture the jump in the Sinusoid-Jump example. We also included another type of step functions, which are the 2d B-spline function with degree equal to zero contained in \mathcal{L}_5 . If a basis selection procedure works well, then it should not select a basis from \mathcal{L}_5 since none of the bases from \mathcal{L}_5 is used to construct the Sinusoid-Jump function. To fit the continuous part of the Sinusoid-Jump

function, we have the Fourier series in \mathcal{L}_2 , the tensor product basis functions in \mathcal{L}_4 , and the cubic spline representers of x_1 in \mathcal{L}_6 . A good basis selection procedure should select a single basis function from \mathcal{L}_4 , and none from \mathcal{L}_6 .

In the simulation, the design points (x_{i1}, x_{i2}) for $i = 1, \dots, n$ are uniformly sampled from $[0, 1] \times [0, 1]$ for the Paraboloid-Jump and Sinusoid-Jump examples. For the Mexican hat example, they are uniformly sampled from $[-4, 4] \times [-4, 4]$. The sample size used is $n \in \{256, 512, 1024\}$. The error variance σ^2 is chosen such that $\text{SNR} \in \{1, 2, 3, 4, 5, 6\}$. For the BSML and LAP procedures, the maximal number of bases M is fixed at 30. For the look-ahead procedure, the IDF spaces are $S_i = \{1, 2, 3\}$ for each library \mathcal{L}_i . M_d is set to be $\prod_{i=1}^5 |S_i| = 3^5 = 243$ for the Paraboloid-Jump and Mexican hat examples, and $\prod_{i=1}^6 |S_i| = 3^6 = 729$ for the Sinusoid-Jump example. The 10-fold cross validation is used to select the final model in the look-ahead procedure.

We also applied the SS-ANOVA model and the MARS procedure to these three examples. For all examples, SS-ANOVA model assumes $f \in W_2^2[0, 1] \otimes W_2^2[0, 1]$, where the support for the Mexican hat example has been scaled to $[0, 1] \times [0, 1]$. For the MARS procedure, we set `degree = 2` in `mars` which means we allow pairwise interactions terms in the model. The IDF used is fixed at 3, and the maximum number of model terms `nk` is fixed at 30. The simulation results for BSML-C, BSML-S, LAP, SS-ANOVA, and MARS based on 100 random samples are shown in Figure 4.24 – 4.26. The average CPU time for each sample size is listed in Table 4.16.

Based on the results, MARS did not work well for all three examples and SS-ANOVA works relatively well only for the Mexican hat example. LAP has the best performance for the Sinusoid-Jump example when $\text{SNR} \geq 3$ across all sample sizes. For the other two examples, LAP has similar performance as the BSML procedures. In terms of computation speed, LAP is faster than BSML-S for all examples, and is similar to BSML-C except for the Sinusoid-Jump example.

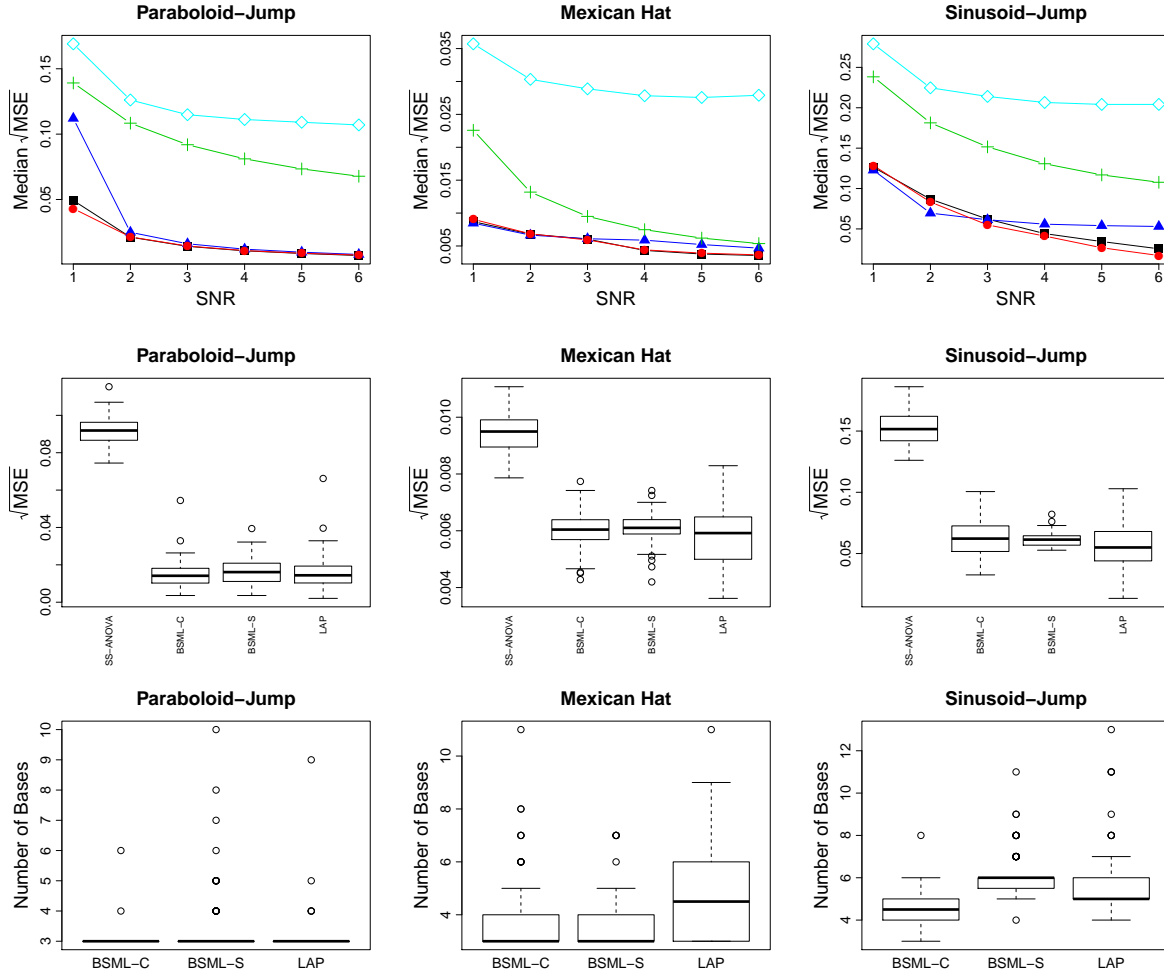


Figure 4.24: Simulation results for the Paraboloid-Jump, Mexican hat, and Sinusoid-Jump examples with $n = 256$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for each example, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, \bullet for LAP, $+$ for SS-ANOVA, and \diamond for MARS. The middle panel shows the boxplots of the $\sqrt{\text{MSE}}$ when SNR = 3. The results of fitting all three examples with MARS are not shown here because they are much worse compared to the results from the BSML and LAP procedures. The bottom panel shows the number of basis functions selected by BSML-C, BSML-S, and LAP when SNR = 3.

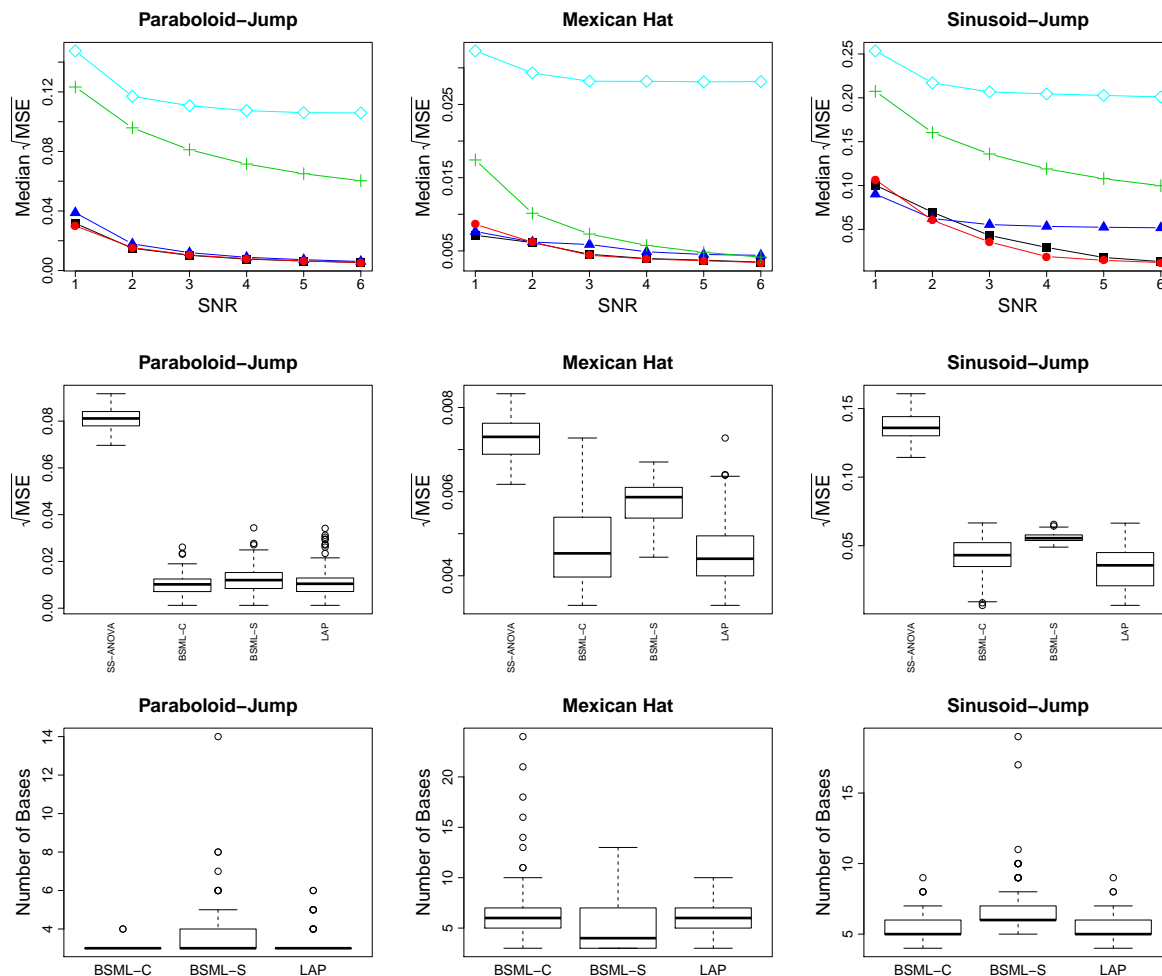


Figure 4.25: Simulation results for the Paraboloid-Jump, Mexican hat, and Sinusoid-Jump examples with $n = 512$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for each example, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, \bullet for LAP, $+$ for SS-ANOVA, and \diamond for MARS. The middle panel shows the boxplots of the $\sqrt{\text{MSE}}$ when SNR = 3. The results of fitting all three examples with MARS are not shown here because they are much worse compared to the results from the BSML and LAP procedures. The bottom panel shows the number of basis functions selected by BSML-C, BSML-S, and LAP when SNR = 3.

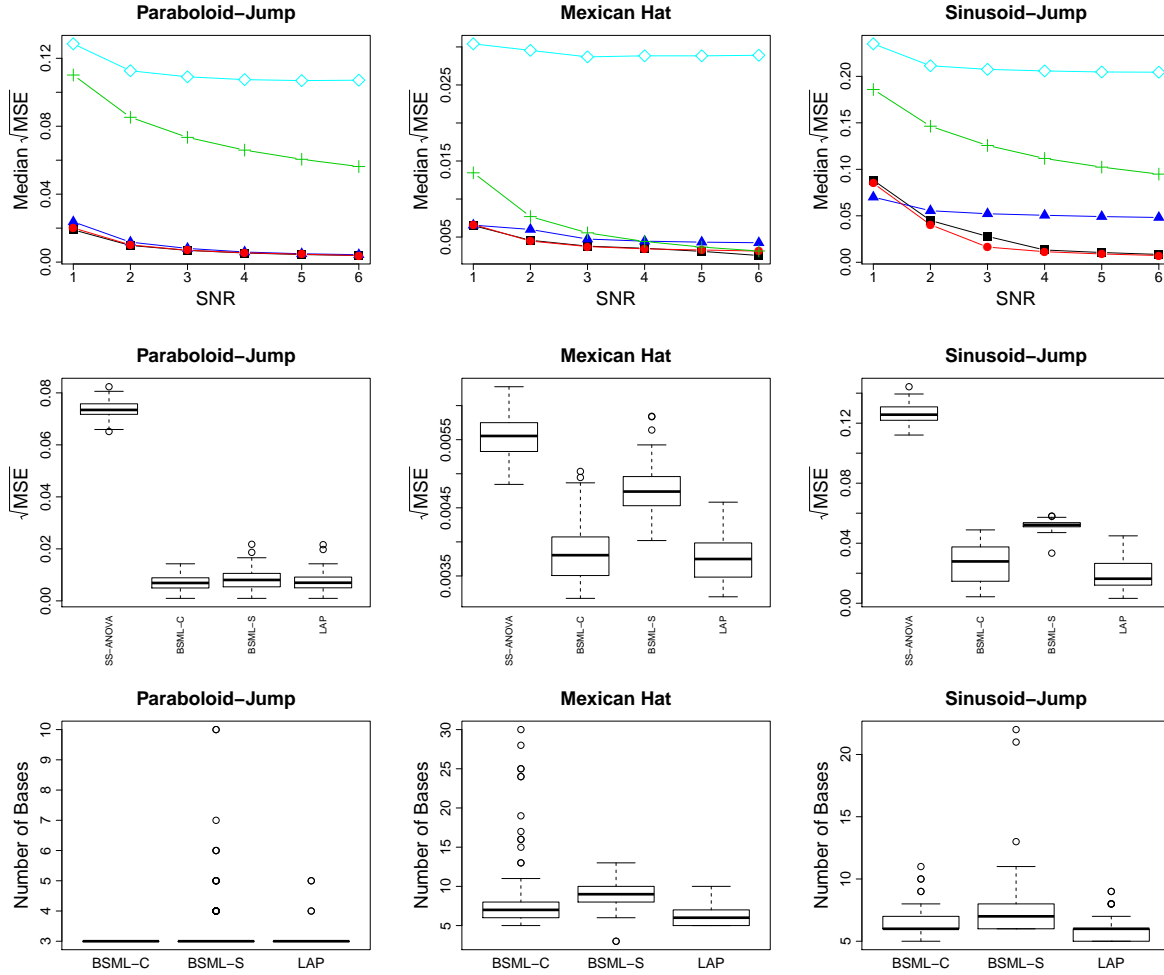


Figure 4.26: Simulation results for the Paraboloid-Jump, Mexican hat, and Sinusoid-Jump examples with $n = 1024$. The top panel shows the median $\sqrt{\text{MSE}}$ across different SNR values for each example, where the symbols used are: \blacksquare for BSML-C, \blacktriangle for BSML-S, \bullet for LAP, $+$ for SS-ANOVA, and \diamond for MARS. The middle panel shows the boxplots of the $\sqrt{\text{MSE}}$ when SNR = 3. The results of fitting all three examples with MARS are not shown here because they are much worse compared to the results from the BSML and LAP procedures. The bottom panel shows the number of basis functions selected by BSML-C, BSML-S, and LAP when SNR = 3.

Table 4.16: Average CPU time (in seconds) of fitting the Paraboloid-Jump, Mexican hat, and Sinusoid-Jump examples using BSML-C, BSML-S, LAP, SS-ANOVA, MARS for sample sizes $n = 256, 512, 1024$ with $\text{SNR} = 3$.

	Paraboloid-Jump			Mexican hat		
	$n = 256$	$n = 512$	$n = 1024$	$n = 256$	$n = 512$	$n = 1024$
BSML-C	1.63	3.06	6.57	1.08	2.12	4.27
BSML-S	15.21	31.62	52.54	13.39	28.64	43.03
LAP	1.40	2.47	4.71	1.55	2.39	3.41
SS-ANOVA	0.85	4.56	36.31	0.47	2.08	9.16
MARS	0.01	0.02	0.03	0.02	0.02	0.04

	Sinusoid-Jump		
	$n = 256$	$n = 512$	$n = 1024$
BSML-C	2.37	4.48	9.71
BSML-S	21.64	43.28	74.05
LAP	6.41	10.13	16.31
SS-ANOVA	0.97	5.50	50.76
MARS	0.01	0.02	0.03

To sum up, based on the simulations presented in Sections 4.1 – 4.4, we conclude that our look-ahead procedure has similar or better performance compared to the BSML procedures, and it tends to select fewer basis functions on average. Moreover, in terms of the computation speed our procedure is superior than BSML-S, and has the edge over BSML-C in some situations.

Chapter 5

Applications

5.1 Well Log

The well log data were introduced in Ruanaidh and Fitzgerald (1996). These data contain measurements of magnetic resonance of the rock strata at different time points while drilling a well. The response magnetic resonance is in the scale of 10^4 . The data set is available at <http://mldata.org/repository/data/viewslug/well-log/>. Because the data contain a considerable number of outliers, we removed the 40 smallest observations as done similarly in Ruanaidh and Fitzgerald (1996). The adjusted data set contains 4010 observations and is plotted in Figure 5.1. The explanatory variable is scaled into $x = [0, 1]$. Here is the R code for loading the data set and creating the variables.

```
library(bsml)
library(assist)
source("LAP.R")

wl <- read.csv("well-log.csv", header=FALSE)
N <- nrow(wl)
outliers <- order(wl[,1])[1:40]
y <- wl[-outliers,1]/(10^4)
```

```

n <- length(y)
time <- c(1:N)[-outliers]
x <- (time-min(time))/(max(time)-min(time))

plot(time, y, type="l", ylab="Magnetic Resonance", xlab="Time Points")

```

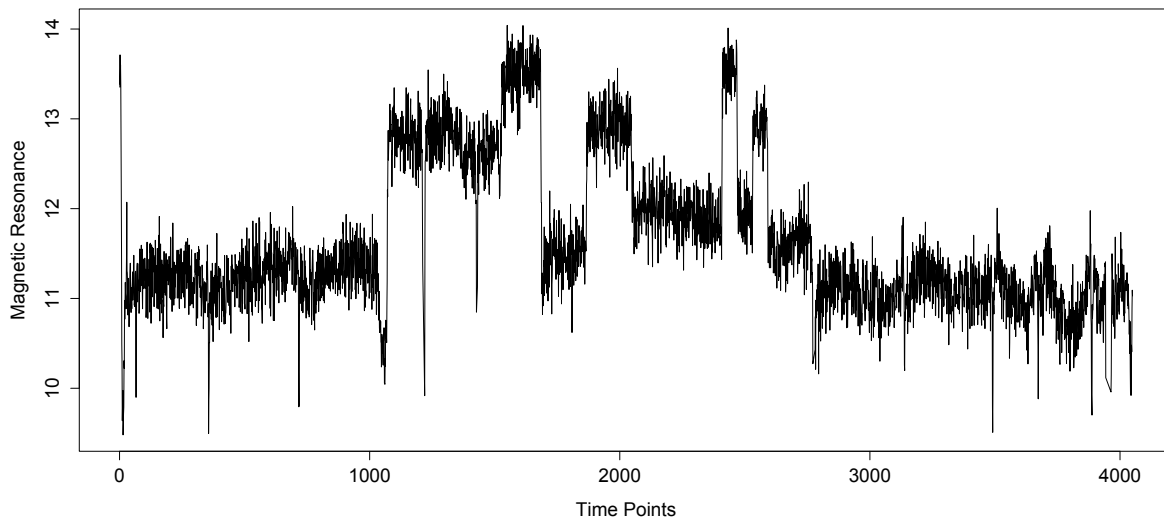


Figure 5.1: Well log data

We first fit a cubic spline to the data using the R function `ssr` in the `assist` package, where the smoothing parameter is selected by the generalized maximum likelihood (GML) criterion. The fit is shown in Figure 5.2. Here is the R code we used:

```

fit.ssr <- ssr(y~x, cubic(x), spar="m")

plot(time, y, col="gray", ylab="Magnetic Resonance", xlab="Time Points",
      main="Cubic Spline", xaxt="n", cex.lab=1.5, cex.main=2)
axis(1, at=seq(0,4100, by=200), label=seq(0,4100, by=200))
lines(time, fit.ssr$fit, lwd=3, col=3)

```

We can see that cubic spline overfits the data. For example, at time points around 1210, and 1430 the fit is clearly affected by some outliers.

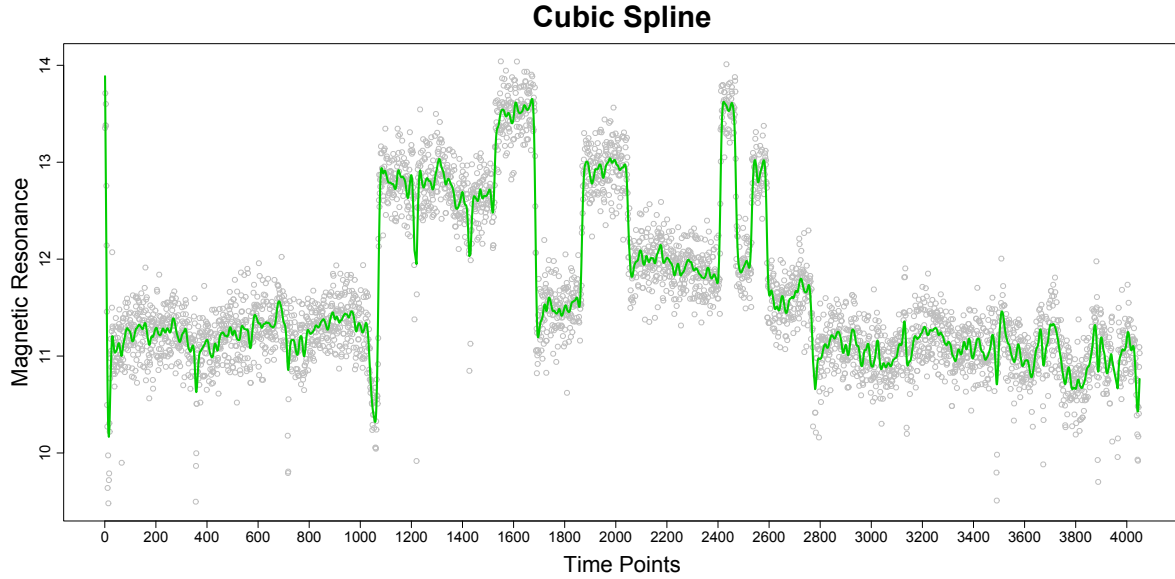


Figure 5.2: Cubic spline fits of the well log data. The Gray dots are the observations. The green line is the cubic spline fit with the smoothing parameter selected by GML criterion.

Next we fit the data using BSML-C, BSML-S and LAP. Table 5.1 lists the libraries of basis functions used in these procedures.

Table 5.1: Basis libraries for the well log example.

$\mathcal{L}_0 = \{1, x\}$
$\mathcal{L}_1 = \mathcal{C}_2(x)$, with knots $\{0.01j, j = 1, \dots, 99\}$
$\mathcal{L}_2 = \mathcal{T}_0(x) = \{1_{(z_1, 1]}(x), \dots, 1_{(z_{99}, 1]}(x)\}$, $z_j = 0.01j$, $j = 1, \dots, 99$
$\mathcal{L}_3 = \mathcal{T}_1(x) = \{(x - z_1)_+, \dots, (x - z_{99})_+\}$, $z_j = 0.01j$, $j = 1, \dots, 99$
$\mathcal{L}_4 = \mathcal{T}_2(x) = \{(x - z_1)_+^2, \dots, (x - z_{99})_+^2\}$, $z_j = 0.01j$, $j = 1, \dots, 99$

Here is how we generate the libraries of basis functions in R:

```
loc <- seq(0.01, 0.99, 0.01)
```

```

baseslist.gen <- function(x, loc){
  n <- length(x)
  pt <- rep(1,n)%o%loc
  L0 <- cbind(1, x)
  L1 <- cubic(x,t=loc)
  L2 <- 1*(x>pt)
  L3 <- (x-pt)*(x>pt)
  L4 <- (x-pt)^2*(x>pt)
  baseslist <- list(L0, L1, L2, L3, L4)
  return(baseslist)
}

baseslist <- baseslist.gen(x, loc)

```

For the BSML and LAP procedures, we specified $M = 100$, which allows at most 100 basis functions to be selected in the final model (including the null bases). Here are the R codes for fitting the data using BSML-C and BSML-S.

```

# BSML-C:
fit.bsmlc <- bsml(y, baseslist, method="bsmlc", maxbas=100)

# BSML-S:
fit.bsmls <- bsml(y, baseslist, method="bsmls", maxbas=100)

# Figure 5.3:
par(mfrow=c(2,1))
plot(time, y, col="gray", ylab="Magnetic Resonance", xlab="Time Points",
      main="BSML-C", xaxt="n", cex.lab=1.5, cex.main=2)
axis(1, at=seq(0,4100, by=200), label=seq(0,4100, by=200))
lines(time, fit.bsmlc$fit, lwd=3, col=1)

plot(time, y, col="gray", ylab="Magnetic Resonance", xlab="Time Points",
      main="BSML-S", xaxt="n", cex.lab=1.5, cex.main=2)
axis(1, at=seq(0,4100, by=200), label=seq(0,4100, by=200))
lines(time, fit.bsmls$fit, lwd=3, col=4)

```

The BSML-C and BSML-S fits are shown in Figure 5.3. The BSML-C procedure also seem to overfit the data. The fit contains many small jumps besides the obvious big jumps. In comparison, the BSML-S procedure gives a smoother fit. BSML-C has

selected 90 basis functions, among which 55 of them are from \mathcal{L}_2 . BSML-S has selected 90 basis functions, among which 20 of them are from \mathcal{L}_2 .

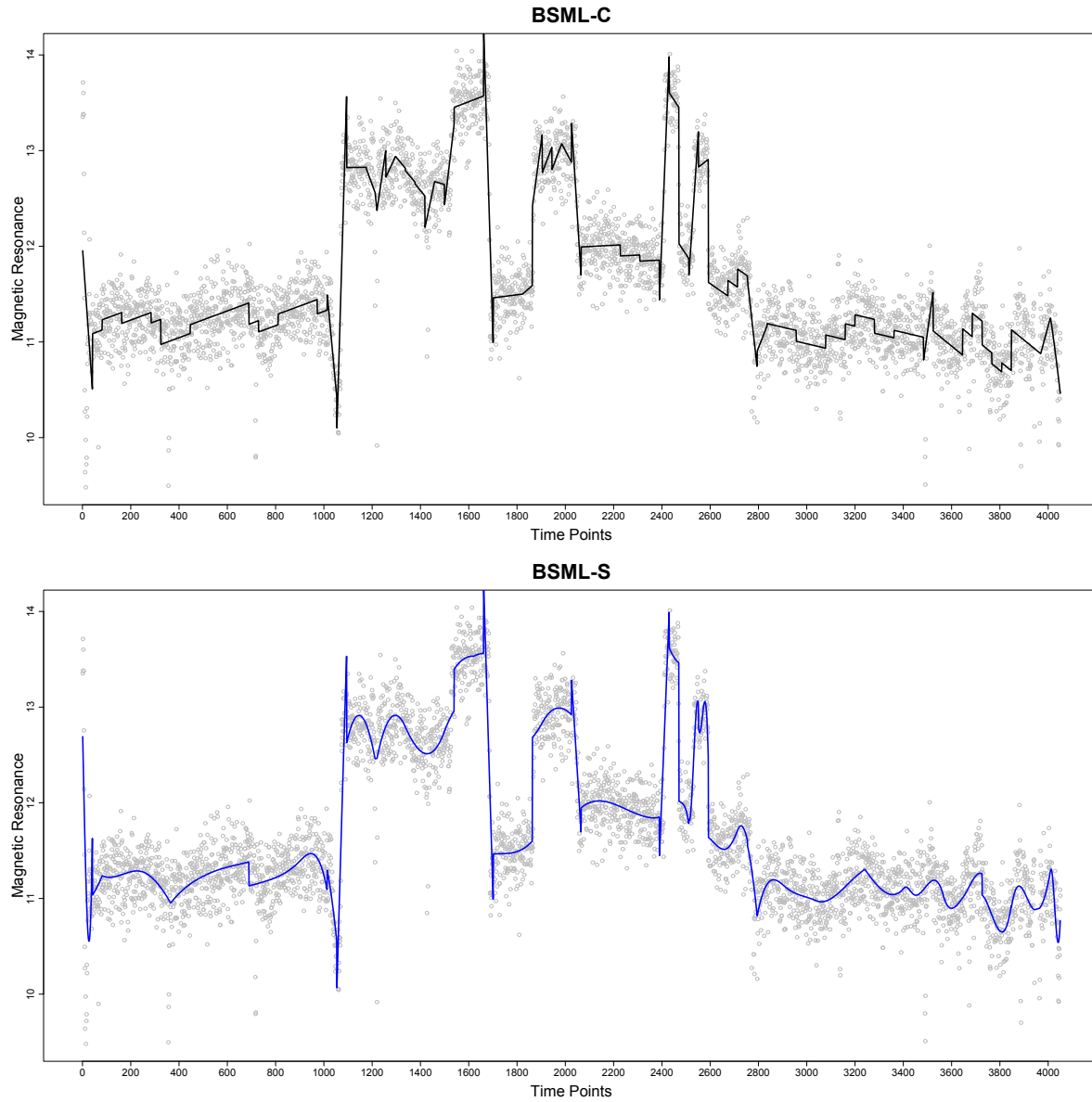


Figure 5.3: BSML-C and BSML-S fits of the well log data. The Gray dots are the observations. The black line in the top panel is the BSML-C fit. The blue line in the lower panel is the BSML-S fit.

For the LAP, the IDF values used are $\{1, 2, 3, 4, 5, 6\}$ for each library, since this works well for the Sine-Jumps and Blocks-Curves examples in our simulation. The input IDF matrix contains in total $6^4 = 1296$ IDF tuning vectors. By the default setting in LAP, $M_d = 1296$, so that given any tuning vector in the IDF matrix, there exists a unique candidate model generated from the forward selection process. Figure 5.4 shows the LAP fit together with the 95% bootstrap confidence interval. Here for the confidence interval we specified the number of bootstrap replicates to be `bootrep = 200`. The R code used is shown below.

```
# LAP:
nlib <- length(baseslist)
idfs <- append(list(1),rep(list(c(1:6)),nlib-1))
idfs <- as.matrix(expand.grid(idfs))

set.seed(123)
fit.lap <- LAP(y, baseslist, maxbas=100, idfs=idfs, criterion="BIC")

# Confidence interval for LAP:
pred.lap <- predict.LAP(fit.lap, bases.include=c(1,1,1,1,1,1),
                        new.baseslist=baseslist, confint=T)

# Figure 5.4:
par(mfrow=c(1,1))
plot(time, y, col="gray", ylab="Magnetic Resonance", xlab="Time Points",
     main="LAP", xaxt="n", cex.lab=1.5, cex.main=2)
axis(1, at=seq(0,4100, by=200), label=seq(0,4100, by=200))
lines(time, pred.lap$fit, lwd=3, col=2)
lines(time, pred.lap$lower, col=1, lwd=2, lty=2)
lines(time, pred.lap$upper, col=1, lwd=2, lty=2)
```

The LAP fit has the most parsimonious representations in terms of the basis functions since only 41 bases are selected, among which 13 of them are from \mathcal{L}_2 . LAP has successfully identified the major jumps and does not overfit the data by adding small jumps like BSML-C does. Although BSML-S does a better job than BSML-C, its fit is still affected by some outliers when compared to LAP, for example at time points around

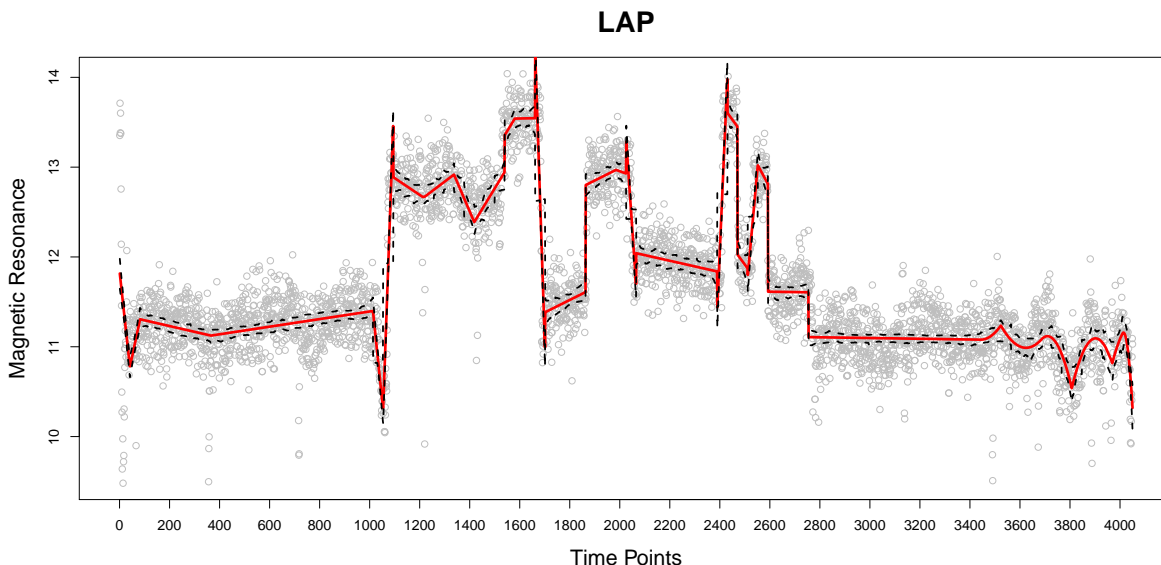


Figure 5.4: LAP fit of the well log data. The Gray dots are the observations. The red solid line is the LAP fit. The black dashed lines are the lower and upper bounds of the 95% bootstrap confidence interval.

40 and 1010.

5.2 Ozone Pressure

The Antarctic ozone hole refers to the ozone depletion observed in the stratosphere of the antarctic region. This phenomenon has been widely studied, for example see Solomon (1999). Ozone depletion in Antarctica was first documented in 1985, and is observed during the Spring season. We will use the ozonesonde data collected at the Amundsen-Scott South Pole Station to investigate this phenomenon. The data is available on the National Oceanic and Atmospheric Administration website: <ftp://ftp.cmdl.noaa.gov/ozwv/Ozonesonde/SouthPole,Antartica/100MeterAverageFiles/>. The raw data between September 2006 and August 2009 contains 0.1-km vertical averages of the ozone partial pressure (in mPa) for 207 unique ozonesonde launch dates. We will focus on the

altitude levels (in km) within the range 11 km to 30 km, where the lower stratosphere is located. Among the total 207 launch dates, we selected 181 dates for which few than 20 values are missing for the ozone partial pressure at the altitude levels of interest.

For our analysis, the response variable is the ozone partial pressure. The covariates are the altitude (from 11 km to 30 km for every 0.1 km, in total 191 levels) and the dates of the records. For convenience, we treat the dates as numerical by counting the number of days elapsed since the start of the records, which is September 06, 2006. We call this variable “day”. For example, for the first record, $\text{day} = 1$. For the second record that occurred on September 09, 2006, $\text{day} = 4$, and so on. The variable “day” has 181 unique values.

Instead of including all the 191 different altitude levels in our analysis, we used two subsets of them. The first subset consists of altitude levels $\{10 + j : j = 1, \dots, 20\}$. The second subset consists of altitude levels $\{10.5 + j : j = 1, \dots, 19\}$. All the observations with altitude levels in the first subset forms the training set. All the observations with altitude levels in the second subset forms the test set. We will build models based on the training set, and test their performances using the test set. For the training set, 91 out of 3620 observations are missing. For the test set, 54 out of 3439 observations are missing. Most of the missing values occur at upper altitudes on a subset of dates. We scale the covariates **altitude** and **day** in both the training set and the test set to $[0, 1]$. Denote x_1 as the scaled **altitude** and x_2 as the scaled **day**.

In the R code below, **alt.train** and **alt.test** are the first and the second subsets of distinct altitude levels. **date.uni** contains unique values of the variable **day**. **ozone.train** is a 20×181 matrix that contains observed ozone partial pressure from the training set. **ozone.test** is a 19×181 matrix that contains observed ozone partial pressure from the test set. The (i, j) th element in **ozone.train** and **ozone.test** corresponds to the ozone pressure observed at the i th altitude on j th ozonesonde launch date. The

missing values in `ozone.train` are removed to form the 3529×1 vector `ozone`. `index.na` contains the indices of these missing values. `x1` and `x2` are 3529×1 vectors that contain the scaled `altitude` and the scaled `day` respectively for each observation in `ozone`. `x1.test` and `x2.test` are the scaled covariates for the observations in `ozone.test`. The code used to generate `x1`, `x2`, `x1.test`, and `x2.test` is given below.

```
library(assist)
library(bsml)
library(gss)
library(lattice)
source("LAP.R")

alt.train <- seq(from=11, to=30, by=1)
n.date <- length(date.uni)
n.alt <- length(alt.train)

index.na <- which(is.na(ozone.train))
ozone <- na.omit(as.vector(ozone.train))
n <- length(ozone)

covar <- as.matrix(expand.grid(alt.train, date.uni))
var1 <- covar[-index.na,1]
var2 <- covar[-index.na,2]

x1 <- (var1-min(var1))/(max(var1)-min(var1))
x2 <- (var2-min(var2))/(max(var2)-min(var2))

alt.test <- seq(from=11.5, to=29.5, by=1)
covar.test <- expand.grid(alt.test, date.uni)
x1.test <- (covar.test[,1]-min(var1))/(max(var1)-min(var1))
x2.test <- (covar.test[,2]-min(var2))/(max(var2)-min(var2))
```

The following R code is used to generate the level plot of the ozone partial pressure, shown in Figure 5.5. Here `month.tick` and `month.labels` are the tick mark locations and labels for the x-axis of the level plot. The level plot illustrates that the lowest ozone partial pressure occurred between 14 km and 20 km in September, October, and November. In addition, the frequency of launches shows seasonality. Among all the launches of interest,

the months with most frequent launches were September and October, with a total of 30 and 31 launches respectively. In contrast, the month with fewest launches was April, which had only 6 launches in total.

Figure 5.5:

```
levelplot(t(ozone.train), aspect=0.5,
  ylab=list(label="Altitude (km)", cex=1.5),
  xlab=list(label="Month/Year", cex=1.5),
  main=list(label="Ozone Pressure from Sep 2006
              to Aug 2009 \n (training set)", cex=2),
  col.regions = topo.colors(300),
  at=seq(min(ozone.train, na.rm=TRUE),
        max(ozone.train, na.rm=TRUE), length=300),
  colorkey=list(tick.number=10, labels=list(cex=1.2)),
  scales=list(x=list(cex=1, at=month.tick, labels=month.labels),
             y=list(cex=1.4, at=1:n.alt))
```

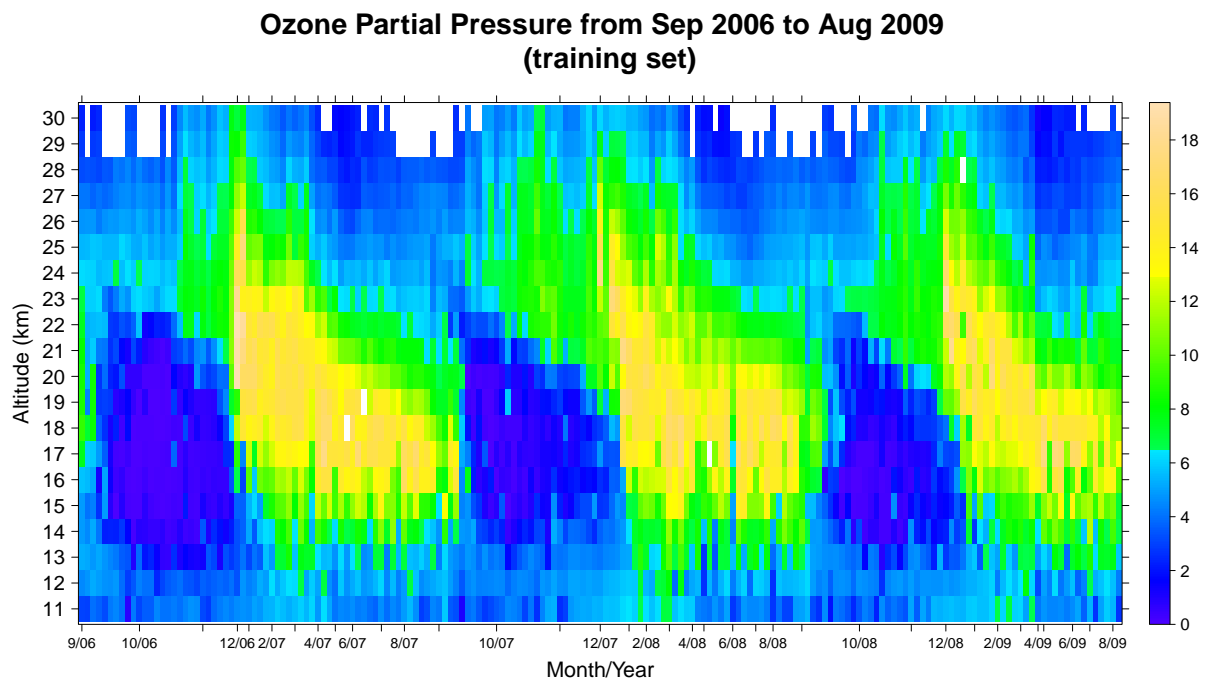


Figure 5.5: Level plot of the ozone partial pressure at altitudes 11 km, 12 km, ..., 30 km, for 181 launch dates between September 2006 and August 2009.

We first fit an SS-ANOVA model with the model space $W_2^2[0, 1] \times W_2^2(per)$ to the training set, using the R package `gss` (Gu (2013a)). The smoothing parameter is selected by the generalized maximum likelihood (GML) criterion. The following R code is used:

```
set.seed(321)
# SS-ANOVA:
z2 <- 3*x2
fit.gss <- ssanova(ozone~x1*z2, type=list(x1=list("cubic", c(0,1)),
                                           z2=list("per", c(0,3))), nbasis=100, method="m")
```

Here we multiply x_2 by 3 to make the period of the periodic spline representers to be $1/3$, so the data for each year are fitted using the same periodic spline representers. In order to make fair comparisons between the performance of the SS-ANOVA model and that of the BSML and LAP procedures, we specified `nbasis=100`, so 100 randomly selected knots are used in the SS-ANOVA model.

Next we fit the training data using the BSML and LAP procedures. Table 5.2 lists the libraries of basis functions used in these procedures. Here \mathcal{L}_1 contains cubic spline representers for the altitude, where the knots used are the unique values of x_1 . \mathcal{L}_2 contains Fourier series with the frequencies equal to multiples of 3, because the data are from three years. The bases in \mathcal{L}_2 can be used to model the seasonal main effect of `day` on the ozone partial pressure. \mathcal{L}_3 contains the basis functions that are tensor products of the bases from $\{x_1 - 0.5\} \cup \mathcal{L}_1$ and the bases from \mathcal{L}_2 . The bases in \mathcal{L}_3 can be used to model the interaction effect of the covariates `altitude` and `day` on the ozone partial pressure. There are 462 basis functions in total.

Table 5.2: Basis libraries for the ozone example.

$\mathcal{L}_0 = \{1, x_1 - 0.5\}$
$\mathcal{L}_1 = \mathcal{C}_2(x_1)$, with knots $\{(j-1)/19, j = 1, \dots, 20\}$
$\mathcal{L}_2 = \{\sin(6\pi k x_2), \cos(6\pi k x_2) : k = 1, 2, \dots, 10\}$
$\mathcal{L}_3 = \{\xi_i \xi_j : \xi_i \in \{x_1 - 0.5\} \cup \mathcal{L}_1, \xi_j \in \mathcal{L}_2\}$

The following code is used to generate the basis libraries for the training set and the test set:

```

baseslist.gen <- function(x1, x2, loc){
  # L0:
  L.null <- cbind(1, x1-0.5)
  # L1:
  L.cubic <- cubic(x1, t=loc)
  # L2:
  kappa <- c(1:10)
  L.fou <- cbind(sin(6*pi*x2%%kappa), cos(6*pi*x2%%kappa))
  # L3:
  L.prod <- (cbind(x1-0.5, L.cubic)%x%matrix(1, ncol=ncol(L.fou)))*
            (matrix(1, ncol=ncol(cbind(x1-0.5, L.cubic))%x%L.fou)
  baseslist <- list(L.null, L.cubic, L.fou, L.prod)
  return(baseslist)
}

# libraries of basis functions for the training set:
baseslist <- baseslist.gen(x1, x2, loc=unique(x1))

# libraries of basis functions for the test set:
baseslist.test <- baseslist.gen(x1.test, x2.test, loc=unique(x1))

```

For the BSML and LAP procedures, we specified the maximum number of bases that can be selected to be $M = 100$. The following code is used to fit the training data with the BSML procedures.

```

set.seed(321)
# BSML-C:

```

```
fit.bsmlc <- bsml(ozone, baseslist, method="bsmlc", maxbas=100)
```

```
# BSML-S:
```

```
fit.bsmls <- bsml(ozone, baseslist, method="bsmls", maxbas=100)
```

For LAP, we first tried the IDF values $\{0.5, 1, 1.5, 2\}$ for each library, so the input IDF matrix contains in total $4^3 = 64$ IDF tuning vectors. By the default setting in LAP, $M_d = 64$, so that each tuning vector in the IDF matrix corresponds to a unique candidate model generated from the forward selection process. The fitted object is named `fit1.lap` in the code below. The IDF values selected by cross validation in `fit1.lap` are 1, 0.5, and 0.5 for \mathcal{L}_1 , \mathcal{L}_2 , and \mathcal{L}_3 respectively. Therefore, we adjusted the IDF values to be $\{0.2, 0.5, 0.8, 1\}$ for each library, and refit the training data with LAP to obtain the final model. The second fitted object is named `fit2.lap`. Given below is the code for fitting the training data twice with LAP.

```
idfs1 <- append(list(1),rep(list(c(0.5, 1, 1.5, 2)),nlib-1))
idfs1 <- as.matrix(expand.grid(idfs1))

set.seed(321)
fit1.lap <- LAP(ozone, baseslist, maxbas=100,
               idfs=idfs1, criterion="BIC")
fit1.lap$idfs.cv

## Adjust the IDF tuning values and refit with LAP:

idfs2 <- append(list(1),rep(list(c(0.2, 0.5, 0.8, 1)),nlib-1))
idfs2 <- as.matrix(expand.grid(idfs2))

set.seed(321)
fit2.lap <- LAP(ozone, baseslist, maxbas=100,
               idfs=idfs2, criterion="BIC")
```

Figures 5.6, 5.7, 5.8, and 5.9 are the level plots of the fitted ozone partial pressure at 191 equally spaced grid points for `altitude` and 218 equally spaced grid points for `day`, by using SS-ANOVA model, BSML-C, BSML-S, and LAP procedures respectively.

Figures 5.10 and 5.11 show the observed and fitted ozone partial pressure at six specific altitudes, so that the differences among SS-ANOVA fit, BSML fits and LAP fit can be seen more easily. The following code is used to generate these figures.

```
## Fitted values at the grid points for SS-ANOVA, BSML-C, BSML-S, LAP:
grid1 <- seq(11, 30, 0.1)
n.grid1 <- length(grid1)
grid2 <- seq(1, 1086, 5)
n.grid2 <- length(grid2)

X.grid <- expand.grid((grid1-11)/(30-11), (grid2-1)/1085)
x1.grid <- X.grid[,1]
x2.grid <- X.grid[,2]

fitted.gss <- matrix(predict(fit.gss, newdata=data.frame(x1=x1.grid,
                                                         z2=3*x2.grid)), nrow=length(grid1), ncol=length(grid2))

baseslist.grid <- baseslist.gen(x1.grid, x2.grid, loc=unique(x1))
fitted.bsmlc <- matrix(predict(fit.bsmlc, bases.include=c(1,1,1,1,1),
                             new.baseslist=baseslist.grid, confint=F),
                       nrow=n.grid1, ncol=n.grid2)
fitted.bsmls <- matrix(predict(fit.bsmls, bases.include=c(1,1,1,1,1),
                             new.baseslist=baseslist.grid, confint=F),
                       nrow=n.grid1, ncol=n.grid2)
fitted.lap <- matrix(predict.LAP(fit2.lap, bases.include=c(1,1,1,1,1),
                               new.baseslist=baseslist.grid, confint=F)$fit,
                     nrow=n.grid1, ncol=n.grid2)

# Tick marks and labels for the plots:
month.tick2 <- c(1, 12, 25, 36, 49, 61, 73, 85, 98,
                110, 122, 134, 146, 159, 171, 183, 195, 207)
month.labels2 <- c("9/06", "11/06", "1/07", "3/07", "5/07", "7/07",
                  "9/07", "11/07", "1/08", "3/08", "5/08", "7/08",
                  "9/08", "11/08", "1/09", "3/09", "5/09", "7/09")
alt.labels <- seq(11,30,1)
alt.tick <- seq(1, 191, 10)

# Figure 5.6:
levelplot(t(fitted.gss), aspect=0.5,
          ylab=list(label="Altitude (km)", cex=1.5),
```

```

xlab=list(label="Month/Year", cex=1.5),
main=list(label="SS-ANOVA fit of the ozone partial
             pressure based on the training set", cex=2),
col.regions = topo.colors(300),
at=seq(min(fitted.gss, na.rm=TRUE),
       max(ozone.train, na.rm=TRUE), length=300),
colorkey=list(tick.number=10, labels=list(cex=1.2)),
scales=list(x=list(cex=1,at=month.tick2, labels=month.labels2),
            y=list(cex=1.4, at=alt.tick, labels=alt.labels)))

```

Figure 5.7:

```

levelplot(t(fitted.bsmlc), aspect=0.5,
          ylab=list(label="Altitude (km)", cex=1.5),
          xlab=list(label="Month/Year", cex=1.5),
          main=list(label="BSML-C fit of the ozone partial
                        pressure based on the training set", cex=2),
          col.regions = topo.colors(300),
          at=seq(min(fitted.bsmlc, na.rm=TRUE),
                 max(ozone.train, na.rm=TRUE), length=300),
          colorkey=list(tick.number=10, labels=list(cex=1.2)),
          scales=list(x=list(cex=1,at=month.tick2, labels=month.labels2),
                      y=list(cex=1.4, at=alt.tick, labels=alt.labels)))

```

Figure 5.8:

```

levelplot(t(fitted.bsmls), aspect=0.5,
          ylab=list(label="Altitude (km)", cex=1.5),
          xlab=list(label="Month/Year", cex=1.5),
          main=list(label="BSML-S fit of the ozone partial
                        pressure based on the training set", cex=2),
          col.regions = topo.colors(300),
          at=seq(min(fitted.bsmls, na.rm=TRUE),
                 max(ozone.train, na.rm=TRUE), length=300),
          colorkey=list(tick.number=10, labels=list(cex=1.2)),
          scales=list(x=list(cex=1,at=month.tick2, labels=month.labels2),
                      y=list(cex=1.4, at=alt.tick, labels=alt.labels)))

```

Figure 5.9:

```

levelplot(t(fitted.lap), aspect=0.5,
          ylab=list(label="Altitude (km)", cex=1.5),

```

```

        xlab=list(label="Month/Year", cex=1.5),
        main=list(label="LAP fit of the ozone partial
                      pressure based on the training set", cex=2),
        col.regions = topo.colors(300),
        at=seq(min(fitted.lap, na.rm=TRUE),
               max(ozone.train, na.rm=TRUE), length=300),
        colorkey=list(tick.number=10, labels=list(cex=1.2)),
        scales=list(x=list(cex=1, at=month.tick2, labels=month.labels2),
                    y=list(cex=1.4, at=alt.tick, labels=alt.labels)))

# Figure 5.10:
par(mfrow=c(3,1))
for(i in c(2,5,8)){
  plot(grid2, fitted.gss[10*(i-1)+1,], type="l",
       col=3, lwd=2, ylim=range(0,17),
       main=paste("Altitude = ", alt.train[i], "km"),
       cex.main=2, ylab="", xlab="", xaxt="n", yaxt="n")
  lines(grid2, fitted.lap[10*(i-1)+1,], col=2, lwd=2)
  lines(grid2, fitted.bsmlc[10*(i-1)+1,], col=1, lwd=2)
  points(date.uni, ozone.train[i,], col="gray60")
  axis(1, at=grid2[month.tick2], labels=month.labels2)
  axis(2, at=seq(0,17,2), labels=seq(0,17,2))
  mtext("Month/Year", side=1, line=3)
  mtext("Ozone Partial Pressure (mPa)", side=2, line=2.5)
}

# Figure 5.11:
par(mfrow=c(3,1))
for(i in c(11,14,17)){
  plot(grid2, fitted.gss[10*(i-1)+1,], type="l",
       col=3, lwd=2, ylim=range(0,17),
       main=paste("Altitude = ", alt.train[i], "km"),
       cex.main=2, ylab="", xlab="", xaxt="n", yaxt="n")
  lines(grid2, fitted.lap[10*(i-1)+1,], col=2, lwd=2)
  lines(grid2, fitted.bsmlc[10*(i-1)+1,], col=1, lwd=2)
  points(date.uni, ozone.train[i,], col="gray60")
  axis(1, at=grid2[month.tick2], labels=month.labels2)
  axis(2, at=seq(0,17,2), labels=seq(0,17,2))
  mtext("Month/Year", side=1, line=3)
  mtext("Ozone Partial Pressure (mPa)", side=2, line=2.5)
}

```

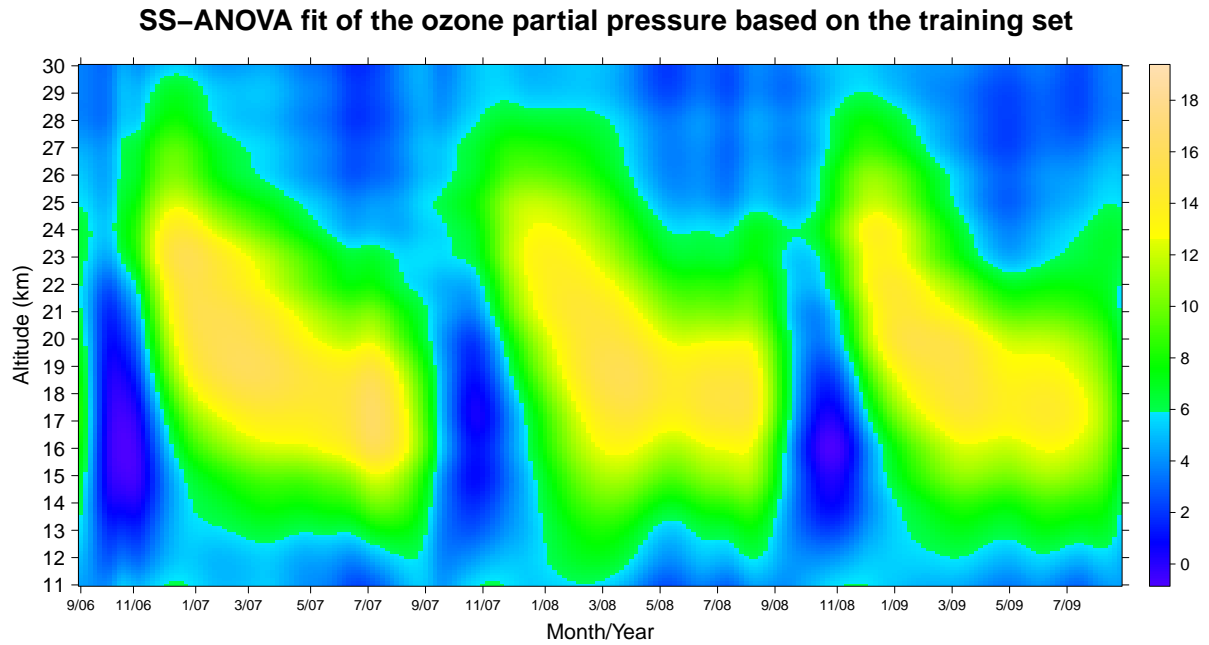


Figure 5.6: Level plot of SS-ANOVA fit of the ozone partial pressure based on the training set, for 191 different altitudes and 218 different dates.

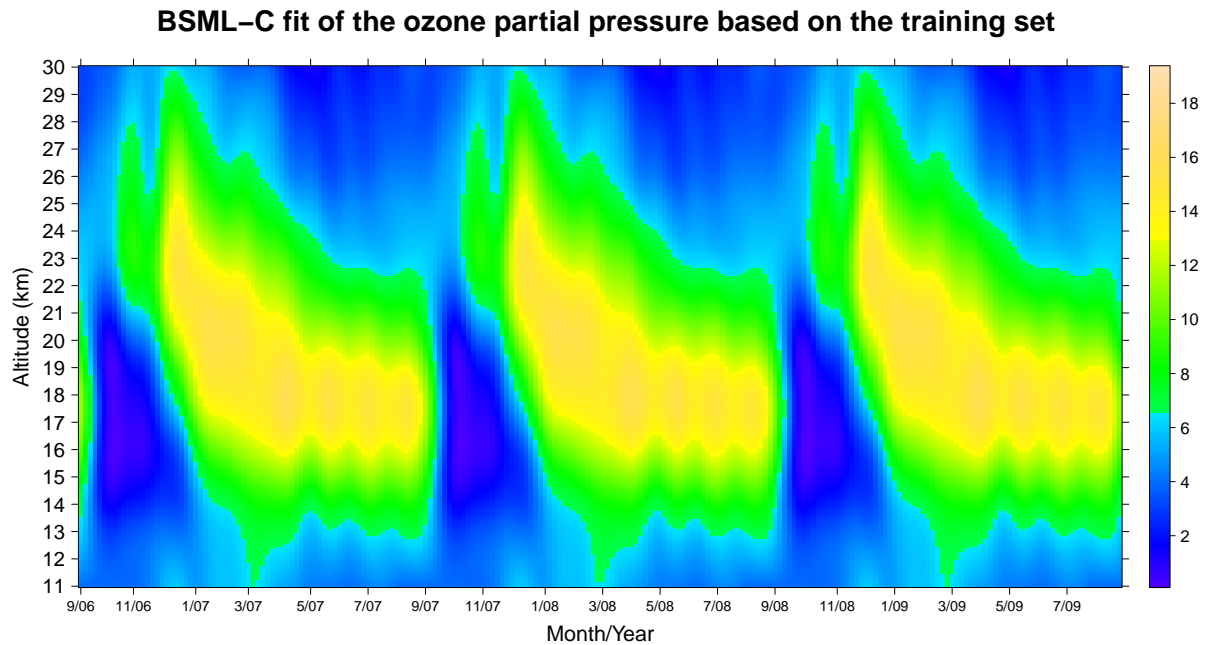


Figure 5.7: Level plot of the BSML-C fit of the ozone partial pressure based on the training set, for 191 different altitudes and 218 different dates.

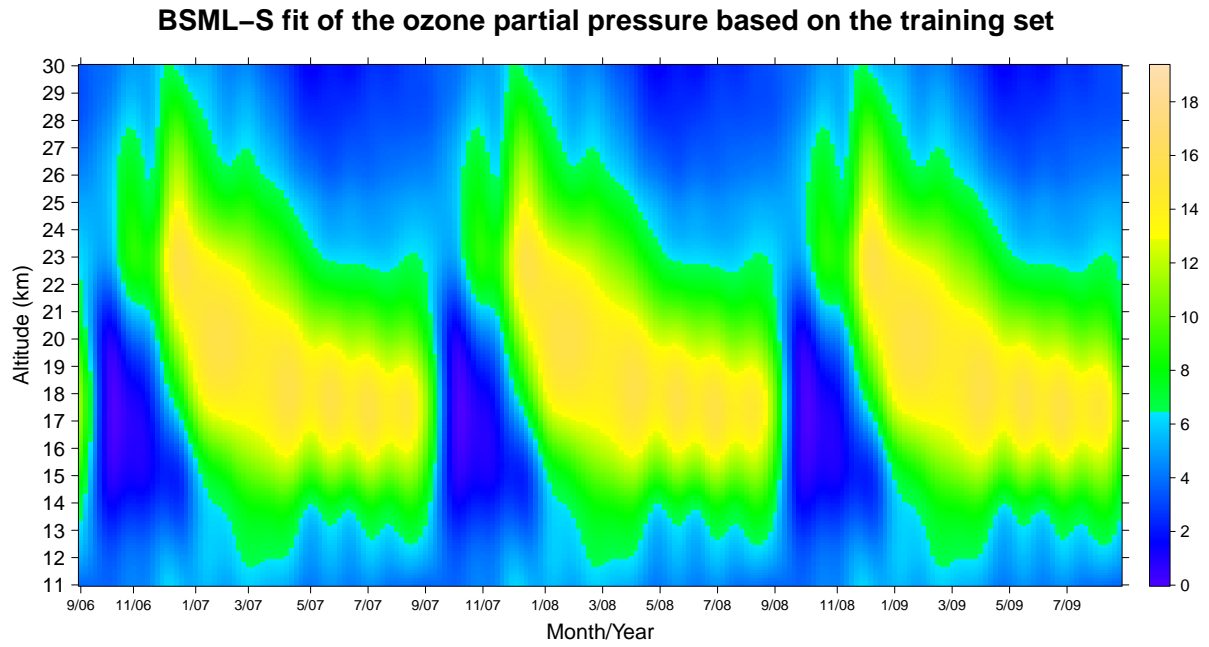


Figure 5.8: Level plot of the BSML-S fit of the ozone partial pressure based on the training set, for 191 different altitudes and 218 different dates.

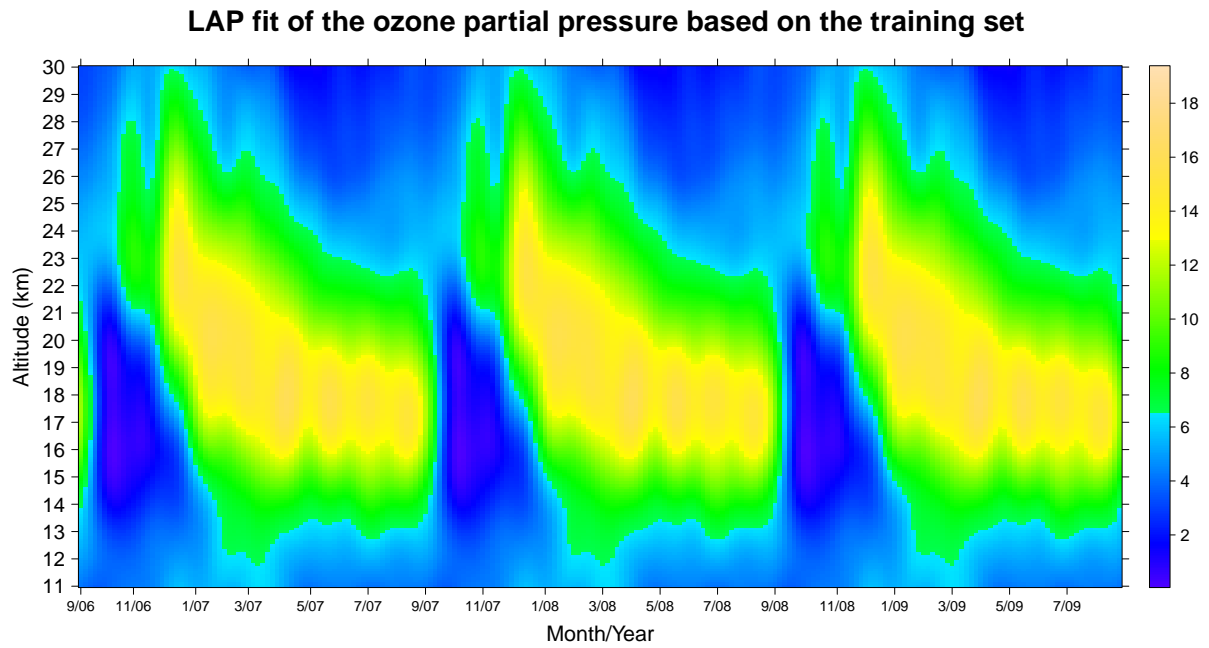


Figure 5.9: Level plot of the LAP fit of the ozone partial pressure based on the training set, for 191 different altitudes and 218 different dates.

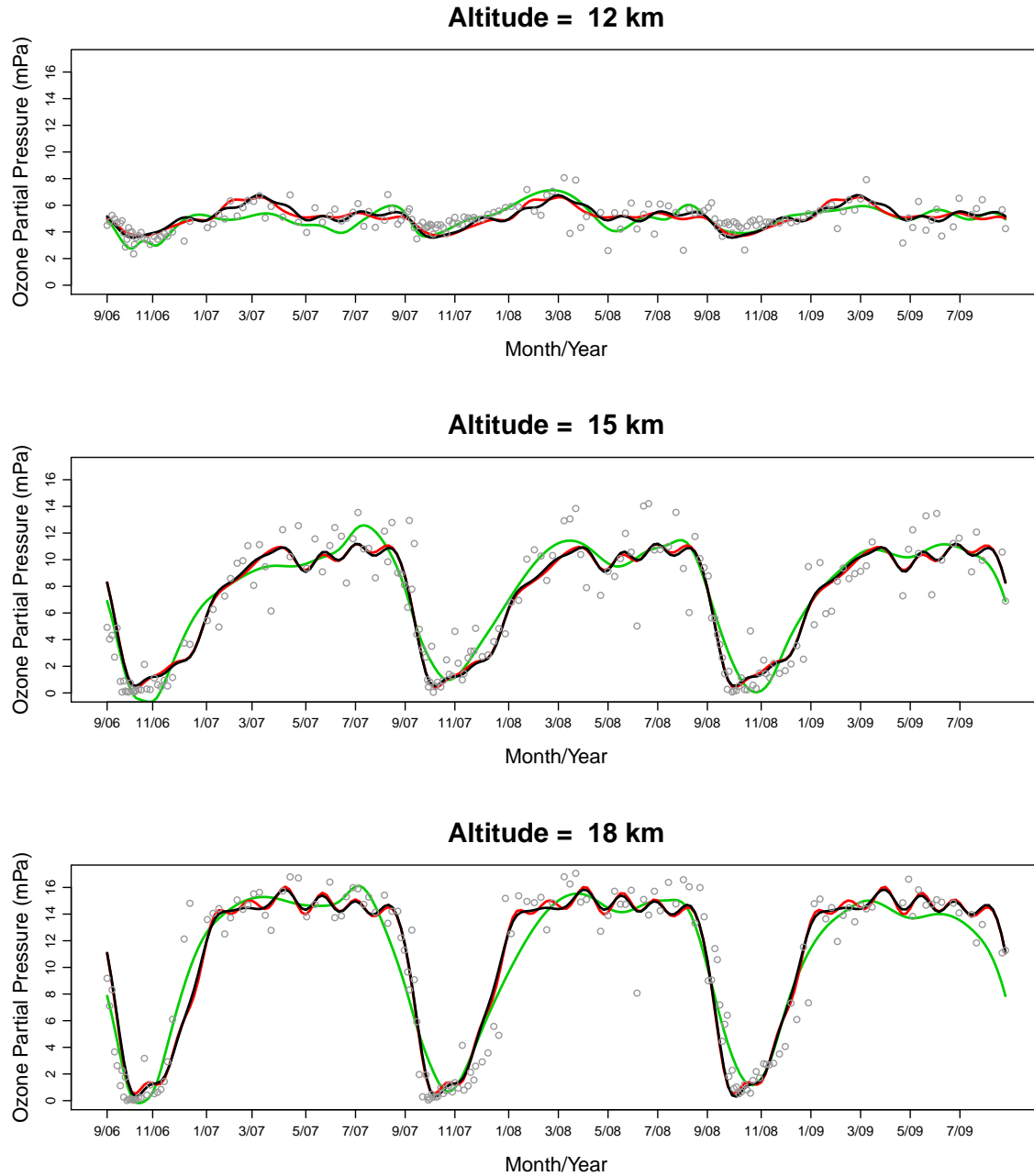


Figure 5.10: SS-ANOVA, BSML-C, BSML-S, and LAP fits of the ozone pressure at altitudes 12 km, 15 km, and 18 km. Here — represents the SS-ANOVA fit, — represents the LAP fit, — represents the BSML-C fit. The BSML-S fit is omitted because it is very close to the BSML-C fit. Circles represent the observations.

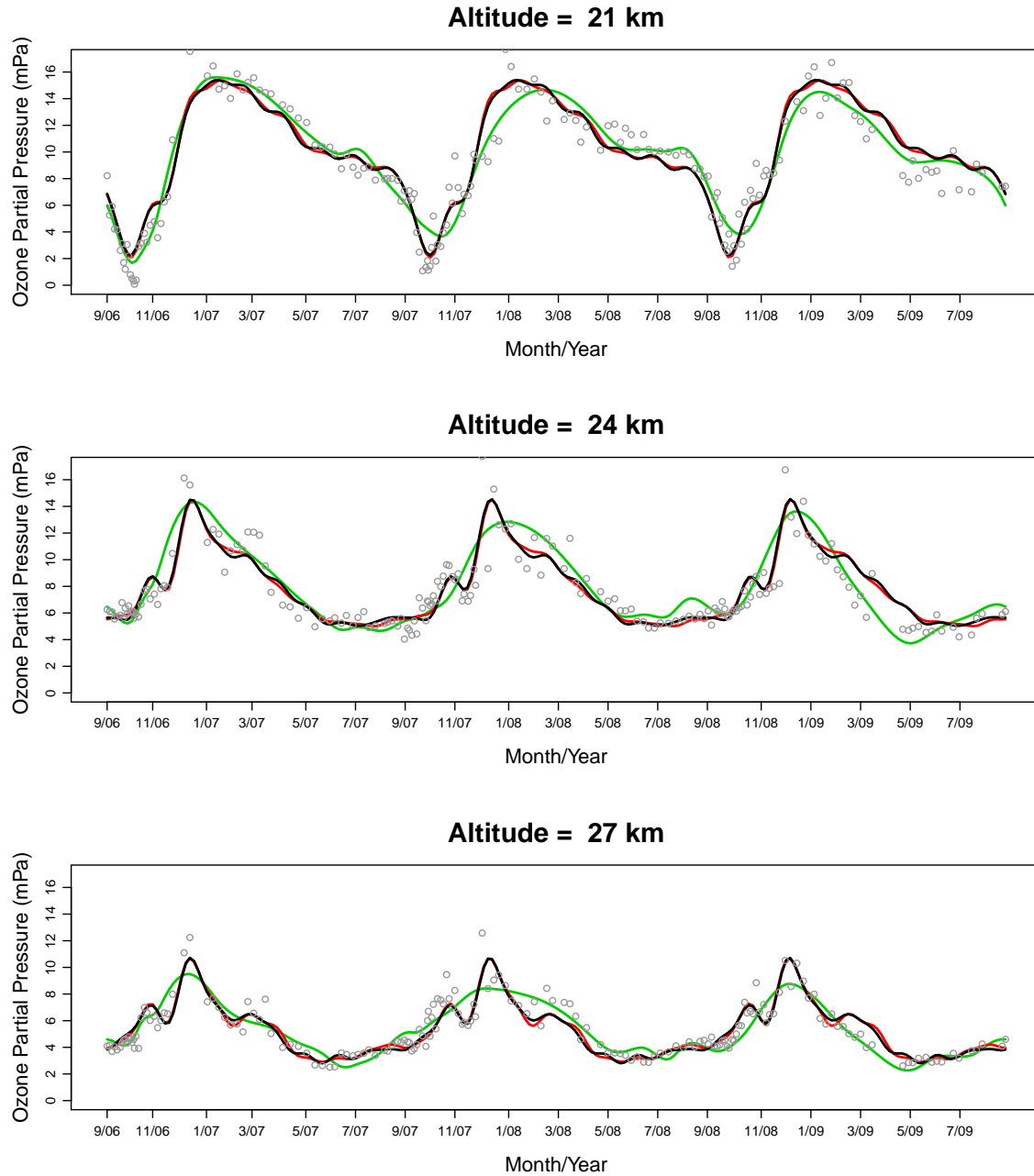


Figure 5.11: SS-ANOVA, BSML-C, BSML-S, and LAP fits of the ozone pressure at altitudes 21 km, 24 km, and 27 km. Here — represents the SS-ANOVA fit, — represents the LAP fit, — represents the BSML-C fit. The BSML-S fit is omitted because it is very close to the BSML-C fit. Circles represent the observations.

Based on the plots above, we can see that the BSML and LAP procedures have similar fits, whereas the SS-ANOVA model tends to oversmooth the data. Among all the 462 basis functions used in the BSML and LAP procedures, 86 of them are selected by LAP. In comparison, BSML-C and BSML-S have selected 67 and 72 bases respectively. The model returned from the LAP procedure has the lowest residual sum of squares compared to other three models. For fair comparison, we compute the predictive squared error (PSE) for the test set using the following code:

```
# PSE from the SS-ANOVA model:
pred.gss <- predict(fit.gss, newdata=data.frame(x1=x1.test,z2=3*x2.test))
( PSE.gss <- mean((pred.gss - ozone.test)^2, na.rm=TRUE) )

# PSE from the BSML-C procedure:
pred.bsmlc <- predict(fit.bsmlc, bases.include=c(1,1,1,1,1),
                     new.baseslist=baseslist.test, confint=F)
( PSE.bsmlc <- mean((pred.bsmlc - as.vector(ozone.test))^2, na.rm=TRUE) )

# PSE from the BSML-S procedure:
pred.bsmls <- predict(fit.bsmls, bases.include=c(1,1,1,1,1),
                     new.baseslist=baseslist.test, confint=F)
( PSE.bsmls <- mean((pred.bsmls - as.vector(ozone.test))^2, na.rm=TRUE) )

# PSE from the look-ahead procedure:
pred.lap <- predict.LAP(fit2.lap, bases.include=c(1,1,1,1,1),
                      new.baseslist=baseslist.test, confint=F)
( PSE.lap <- mean((pred.lap$fit - as.vector(ozone.test))^2, na.rm=TRUE) )
```

We found that $\text{PSE.gss} = 2.191$, $\text{PSE.bsmlc} = 1.863$, $\text{PSE.bsmls} = 1.900$, and $\text{PSE.lap} = 1.826$. Therefore, because the LAP procedure returns the lowest PSE, we conclude that LAP has the best performance compared to other three methods.

For the look-ahead procedure, here we construct the 95% bootstrap confidence intervals for the seasonal and nonseasonal effects at the same grid points of `altitude` and `day` used earlier. The seasonal effect includes the main effect of `day` and the interaction effect between `altitude` and `day`, which can be obtained by using the bases in \mathcal{L}_2 and \mathcal{L}_3 that

have been selected by LAP. We can estimate the seasonal effect alone by specifying the last two elements in the `bases.include` argument of the `predict.LAP` function to be 1 and others be 0, as shown in the code below. The nonseasonal effect includes the constant and the main effect of the latitude, which can be obtained by using the bases in \mathcal{L}_0 and \mathcal{L}_1 that have been selected by LAP. We can estimate the nonseasonal effect alone by specifying the first three elements in the `bases.include` argument of the `predict.LAP` function to be 1 and others be 0, where the first two elements in the `bases.include` corresponds to the null basis functions. For the same six altitudes studied earlier, the seasonal effects are plotted in Figure 5.12 and 5.13. The nonseasonal effect is shown in 5.14. Here is the R code we used to generate these figures.

```
# Seasonal fit and its confidence interval:
lap.season.fit <- predict.LAP(fit2.lap, bases.include=c(0,0,0,1,1),
new.baseslist=baseslist.grid, confint=T)
season.lap <- matrix(lap.season.fit$fit, nrow=n.grid1, ncol=n.grid2)
season.lap.lower <- matrix(lap.season.fit$lower,
                           nrow=n.grid1, ncol=n.grid2)
season.lap.upper <- matrix(lap.season.fit$upper,
                           nrow=n.grid1, ncol=n.grid2)

# Nonseasonal fit and its confidence interval:
lap.nonseason.fit <- predict.LAP(fit2.lap, bases.include=c(1,1,1,0,0),
new.baseslist=baseslist.grid, confint=T)
nonseason.lap <- matrix(lap.nonseason.fit$fit, nrow=n.grid1, ncol=n.grid2)
nonseason.lap.lower <- matrix(lap.nonseason.fit$lower,
                              nrow=n.grid1, ncol=n.grid2)
nonseason.lap.upper <- matrix(lap.nonseason.fit$upper,
                              nrow=n.grid1, ncol=n.grid2)

# Figure 5.12:
par(mfrow=c(3,1))
for(i in c(11, 41, 71)){
  plot(grid2, season.lap[i,], col=2, lwd=2,
       main=paste("Seasonal Effect \n Altitude = ", grid1[i], "km"),
       ylim=range(min(season.lap.lower[i,], na.rm=T),
                  max(season.lap.upper[i,], na.rm=T)),
       xlab="", ylab="", cex.main=1.5, type="l", xaxt="n")
}
```

```

    axis(1,at=grid2[month.tick2], labels=month.labels2)
    lines(grid2, season.lap.lower[i,], col=1, lwd=2, lty=2)
    lines(grid2, season.lap.upper[i,], col=1, lwd=2, lty=2)
    mtext("Month/Year", side=1, line=3)
    mtext("Ozone Partial Pressure (mPa)", side=2, line=2.5)
}

# Figure 5.13:
par(mfrow=c(3,1))
for(i in c(101,131,161)){
  plot(grid2, season.lap[i,], col=2, lwd=2,
       main=paste("Seasonal Effect \n Altitude = ", grid1[i], "km"),
       ylim=range(min(season.lap.lower[i,], na.rm=T),
                  max(season.lap.upper[i,], na.rm=T)),
       xlab="", ylab="", cex.main=1.5, type="l", xaxt="n")
  axis(1,at=grid2[month.tick2], labels=month.labels2)
  lines(grid2, season.lap.lower[i,], col=1, lwd=2, lty=2)
  lines(grid2, season.lap.upper[i,], col=1, lwd=2, lty=2)
  mtext("Month/Year", side=1, line=3)
  mtext("Ozone Partial Pressure (mPa)", side=2, line=2.5)
}

# Figure 5.14:
par(mfrow=c(1,1))
plot(nonseason.lap[,1], grid1, col=2, lwd=2, main="Nonseasonal Effect",
     xlim=range(min(nonseason.lap.lower[,1]-1, na.rm=T),
                max(nonseason.lap.upper[,1]+1, na.rm=T)),
     xlab="", ylab="", cex.main=1.5, type="l", yaxt="n", yaxt="n")
axis(1, at=2:12, labels=2:12)
axis(2, at=alt.train, labels=alt.train)
lines(nonseason.lap.lower[,1], grid1, col=1, lwd=2, lty=2)
lines(nonseason.lap.upper[,1], grid1, col=1, lwd=2, lty=2)
mtext("Altitude (km)", side=2, line=3)
mtext("Ozone Partial Pressure (mPa)", side=1, line=2.5)

```

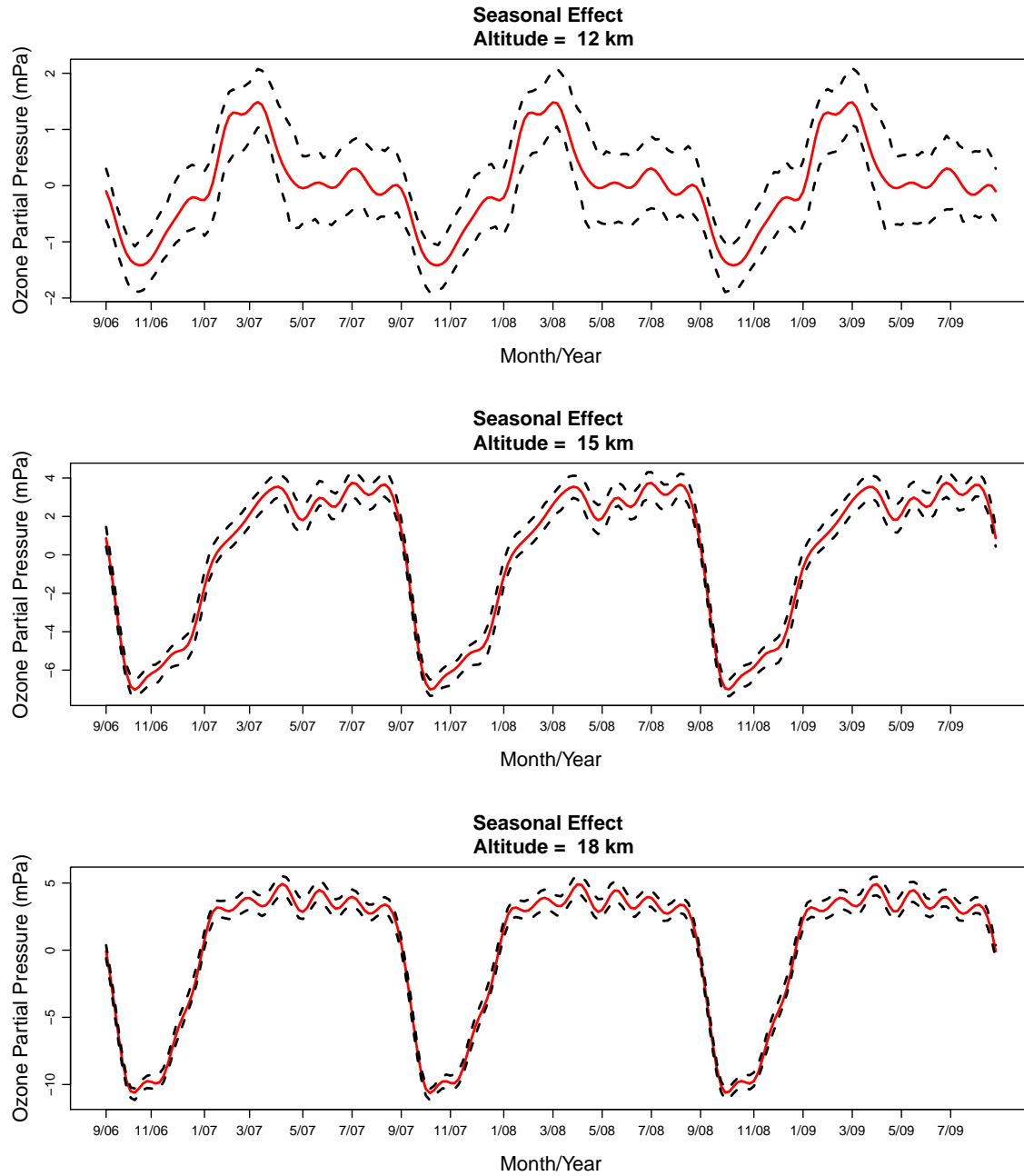


Figure 5.12: Seasonal effect at altitudes 12 km, 15 km, 18 km. The red solid line is the fit. The black dashed lines are the lower and upper bounds of the 95% bootstrap confidence interval.

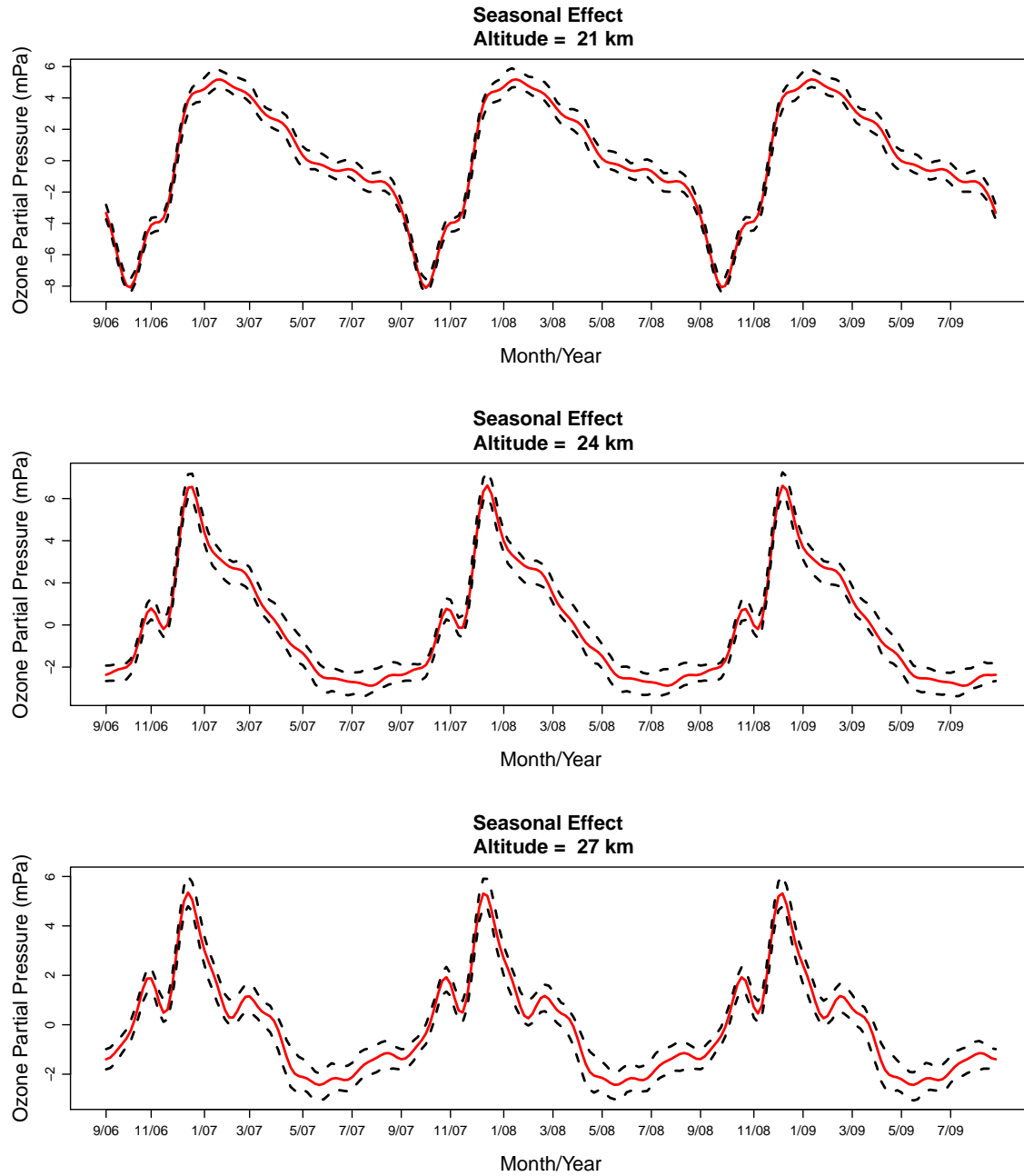


Figure 5.13: Seasonal effect at altitudes 21 km, 24 km, and 27 km. The red solid line is the fit. The black dashed lines are the lower and upper bounds of the 95% bootstrap confidence interval.

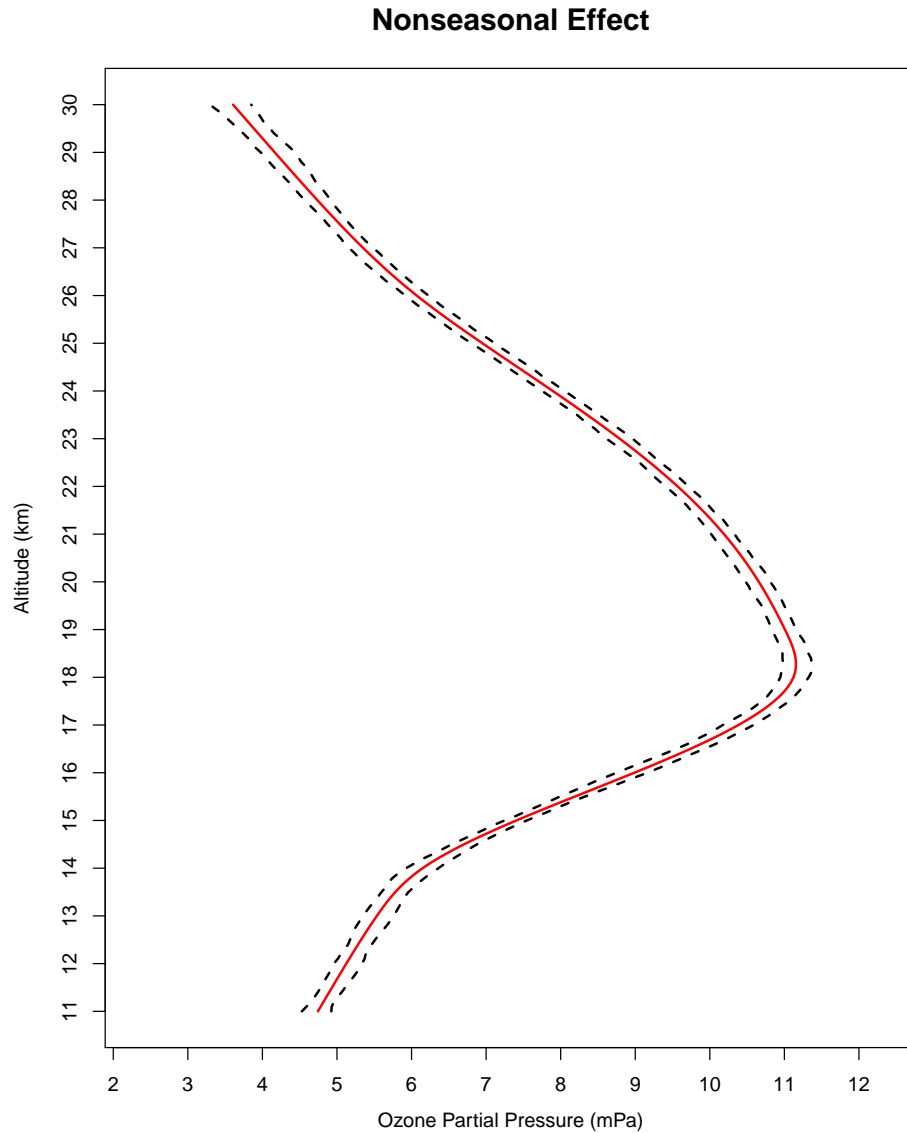


Figure 5.14: Nonseasonal effect. The red solid line is the fit. The black dashed lines are the lower and upper bounds of the 95% bootstrap confidence interval.

Based on Figures 5.12, 5.13, and 5.14, we can see that the seasonal effect is different at different altitudes. From Figure 5.12 we can observe the clear drop of ozone pressure during Spring time at low altitudes. However, this pattern no longer exists at high altitudes, as shown in the middle and bottom panels of Figure 5.13. The nonseasonal

effect does not depend on the dates, since it is fitted using only basis functions related to the variable `altitude`. Figure 5.14 suggests that in general, we should expect the ozone partial pressure to first increase and then decrease as the altitude increases. This finding agrees with the general pattern of ozone distribution in the stratosphere reported in literature. For example, Watson et al. (1986) mentioned that “the ozone layer has a continuous distribution with a peak concentration in the lower stratosphere between about 20 and 25 kilometers altitude”.

Bibliography

- H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6):716–723, 1974.
- D. P. Bertsekas, J. N. Tsitsiklis, and C. Wu. Rollout algorithms for combinatorial optimization. *Journal of Heuristics*, 3(3):245–262, 1997. doi: <http://dx.doi.org/10.1023/A:1009635226865>.
- S. Bjerve, K. A. Doksum, and B. S. Yandell. Uniform confidence bounds for regression based on a simple moving average. *Scandinavian Journal of Statistics*, 12(2):159–169, 1985.
- G. Claeskens and I. V. Keilegom. Bootstrap confidence bands for regression curves and their derivatives. *The Annals of Statistics*, 31(6):1852–1884, 2003.
- P. Craven and G. Wahba. Smoothing noisy data with spline functions. *Numerische Mathematik*, 31(4):377–403, 1978.
- D. J. Cummins, T. G. Filloon, and D. Nychka. Confidence intervals for nonparametric curve estimates: toward more uniform pointwise coverage. *Journal of the American Statistical Association*, 96(453):233–246, 2001.
- B. Efron. Nonparametric standard errors and confidence intervals. *Canadian Journal of Statistics*, 9(2):139–158, 1981. doi: <http://dx.doi.org/10.2307/3314608>.
- B. Efron. *The Jackknife, the Bootstrap and Other Resampling Plans*, volume CBMS-NSF Regional Conference Series in Applied Mathematics, Vol. 38. SIAM, Philadelphia, 1982.
- R. L. Eubank. *Spline Smoothing and Nonparametric Regression*. Marcel Dekker, second edition, 1999.
- R. L. Eubank and P. L. Speckman. Confidence bands in nonparametric regression. *Journal of the American Statistical Association*, 88(424):1287–1301, 1993.
- J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001.

- J. Friedman. Multivariate adaptive regression splines. *The Annals of Statistics*, 19(1): 1–67, 1991.
- J. Friedman, T. Hastie, H. Höfling, and R. Tibshirani. Pathwise coordinate optimization. *The Annals of Applied Statistics*, 1(2):302–332, 2007. doi: <http://dx.doi.org/10.1214/07-AOAS131>.
- J. Friedman, T. Hastie, and R. Tibshirani. *glmnet: Regularization Paths for Generalized Linear Models via Coordinate Descent*, 2010. URL <http://www.jstatsoft.org/v33/i01/>. R package version 2.0-5.
- D. Frost and R. Dechter. Look-ahead value ordering for constraint satisfaction problems. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*, pages 572–578, 1995.
- T. Gasser, L. Sroka, and C. Jennen-Steinmetz. Residual variance and residual pattern in nonlinear regression. *Biometrika*, 73(3):625–633, 1986.
- G. H. Golub and C. F. V. Loan. *Matrix Computations*. The Johns Hopkins University Press, fourth edition, 2013.
- C. Gu. *gss: General Smoothing Splines*, 2013a. URL <http://CRAN.R-project.org/package=gss>. R package version 2.1-0.
- C. Gu. *Smoothing Spline ANOVA Models*. Springer-Verlag New York, second edition, 2013b.
- P. Hall. Effect of bias estimation on coverage accuracy of bootstrap confidence intervals for a probability density. *The Annals of Statistics*, 20(2):675–694, 1992a.
- P. Hall. On bootstrap confidence intervals in nonparametric regression. *The Annals of Statistics*, 20(2):695–711, 1992b.
- P. Hall and J. Horowitz. A simple bootstrap method for constructing nonparametric confidence bands for functions. *The Annals of Statistics*, 41(1):1892–1921, 2013.
- W. Härdle and A. W. Bowman. Bootstrapping in nonparametric regression: local adaptive smoothing and confidence bands. *Journal of the American Statistical Association*, 83(401):102–110, 1988.
- W. Härdle and J. S. Marron. Bootstrap simultaneous error bars for nonparametric regression. *The Annals of Statistics*, 19(2):778–796, 1991.
- W. Härdle, S. Huet, and E. Jolivet. Better bootstrap confidence intervals for regression curve estimation. *Statistics*, 26(4):287–306, 1995.

- T. Hastie and J. Qian. Glmnet vignette, 2014. URL http://www.stanford.edu/~hastie/glmnet/glmnet_alpha.html.
- T. Hastie and R. Tibshirani. *Generalized Additive Models*. Chapman and Hall, 1990.
- T. Hastie, R. Tibshirani, F. Leisch, K. Hornik, and B. D. Ripley. *mda: Mixture and flexible discriminant analysis*, 2015. URL <http://CRAN.R-project.org/package=mda>. S original by Trevor Hastie and Robert Tibshirani. Original R port by Friedrich Leisch and Kurt Hornik and Brian D. Ripley. R package version 0.4-8.
- Y. Lin and H. H. Zhang. Component selection and smoothing in multivariate nonparametric regression. *The Annals of Statistics*, 34(5):2272–2297, 2006. doi: <http://dx.doi.org/10.1214/009053606000000722>.
- Z. Luo and G. Wahba. Hybrid adaptive splines. *Journal of the American Statistical Association*, 92(437):107–116, 1997.
- T. L. McMurphy and D. N. Politis. Bootstrap confidence intervals in nonparametric regression with built-in bias correction. *Statistics & Probability Letters*, 78(15):2463–2469, 2008.
- M. H. Neumann. Automatic bandwidth choice and confidence intervals in nonparametric regression. *The Annals of Statistics*, 23(6):1937–1959, 1995.
- M. H. Neumann and J. Polzehl. Simultaneous bootstrap confidence bands in nonparametric regression. *Journal of Nonparametric Statistics*, 9(4):307–333, 1998.
- D. Picard and K. Tribouley. Adaptive confidence interval for pointwise curve estimation. *The Annals of Statistics*, 28(1):298–335, 2000.
- W. Qian and Y. Yang. Model selection via standard error adjusted adaptive lasso. *Annals of the Institute of Statistical Mathematics*, 65(2):295–318, 2013.
- P. Qiu. *Image Processing and Jump Regression Analysis*. John Wiley & Sons Publication, 2005.
- J. Rice. Bandwidth choice for nonparametric regression. *The Annals of Statistics*, 12(4):1215–1230, 1984.
- J. Ruanaidh and W. Fitzgerald. *Numerical Bayesian Methods Applied to Signal Processing*. Springer-Verlag New York, 1996.
- G. Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464, 1978.
- G. A. Seber and A. J. Lee. *Linear Regression Analysis*. John Wiley & Sons Publication, second edition, 2003.

- J. C. Sklar, J. Wu, W. Meiring, and Y. Wang. Non-parametric regression with basis selection from multiple libraries. *Technometrics*, 55(2):189–201, 2013. doi: <http://dx.doi.org/10.1080/00401706.2012.739104>.
- S. Solomon. Stratospheric ozone depletion: a review of concepts and history. *Reviews of Geophysics*, 37(3):275–316, 1999.
- J. Sun and C. R. Loader. Simultaneous confidence bands for linear regression and smoothing. *The Annals of Statistics*, 22(3):1328–1345, 1994.
- T. Tong and Y. Wang. Estimating residual variance in nonparametric regression using least squares. *Biometrika*, 92(4):821–830, 2005.
- S. Voss, A. Fink, and C. Duin. Looking ahead with the pilot method. *Annals of Operations Research*, 136(1):285–302, 2005. doi: <http://dx.doi.org/10.1007/s10479-005-2060-2>.
- G. Wahba. Bayesian confidence intervals for the cross-validated smoothing spline. *Journal of the Royal Statistical Society, Series B*, 45(1):133–150, 1983.
- G. Wahba. *Spline models for observational data*, volume CBMS-NSF Regional Conference Series in Applied Mathematics, Vol. 59. SIAM, Philadelphia, 1990.
- Y. Wang. *Smoothing Splines Methods and Applications*. CRC Press, 2011.
- Y. Wang and C. Ke. *assist: A Suite of S-Plus Functions Implementing Smoothing Splines*, 2015. URL <http://CRAN.R-project.org/package=assist>. R package version 3.1.3.
- Y. Wang and G. Wahba. Bootstrap confidence intervals for smoothing splines and their comparison to bayesian confidence intervals. *Journal of Statistical Computation and Simulation*, 51(2-4):263–279, 1995.
- R. T. Watson, M. A. Geller, R. S. Stolarski, and R. F. Hampson. Present state of knowledge of the upper atmosphere: An assessment report. *NASA Reference Publication 1162*, 1986.
- J. Wu, J. Sklar, Y. Wang, and W. Meiring. *bsml: Basis Selection from Multiple Libraries*, 2012. URL <https://cran.r-project.org/src/contrib/Archive/bsml/>. R package version 1.5-1.
- Y. Xia. Bias-corrected confidence bands in nonparametric regression. *Journal of the Royal Statistical Society. Series B. Statistical Methodology*, 60(4):797–811, 1998.
- J. Ye. On measuring and correcting the effects of data mining and model selection. *Journal of the American Statistical Association*, 93(441):120–131, 1998.

- H. H. Zhang and C.-Y. Lin. *cosso: Fit Regularized Nonparametric Regression Models Using COSSO Penalty*, 2013. URL <http://CRAN.R-project.org/package=cosso>. R package version 2.1-1.
- J. L. Zhang and J. S. Liu. A new sequential importance sampling method and its application to the two-dimensional hydrophobic-hydrophilic model. *Journal of Chemical Physics*, 117(7):3492–3498, 2002.
- J. L. Zhang, M. T. Lin, J. S. Liu, and R. Chen. Lookahead and piloting strategies for variable selection. *Statistica Sinica*, 17(3):985–1003, 2007.
- H. Zou. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101(476):1418–1429, 2006.